

Lucrările conferinței
*Resurse lingvistice și instrumente pentru
prelucrarea limbii române*
București, 6–7 mai 2010

Volum apărut cu sprijinul Ministerului Educației, Cercetării și
Tineretului, prin Centrul Național de Management Programe
(CNMP), proiectul eDTLR – Dicționarul Tezaur al Limbii
Române în format electronic (contract 91–013/18.09.2007)

ISSN 1843-911X

Lucrările conferinței
*Resurse lingvistice și instrumente pentru
prelucrarea limbii române*
București, 6-7 mai 2010

Editori:
Adrian Iftene
Horia-Nicolai Teodorescu
Dan Cristea
Dan Tufiș

Organizatori:
Facultatea de Informatică,
Universitatea „Alexandru Ioan Cuza” Iași

Institutul de Cercetări pentru Inteligență Artificială
Academia Română, București

Institutul de Informatică Teoretică
Academia Română, Filiala Iași

Muzeul Național al Literaturii Române
Filiala București

Intelligentics, Cluj-Napoca

Editura Universității „Alexandru Ioan Cuza” Iași

COMITETUL DE PROGRAM:

Vasile Apopei, Institutul de Informatică Teoretică, A.R., Iași
Verginica Barbu Mititelu, Institutul de Cercetări în Inteligență Artificială, A.R., București
Corneliu Burileanu, Facultatea de Electronică, Universitatea Politehnică București și Institutul de Cercetări în Inteligență Artificială, A.R., București
Monica Busuioc, Institutul de Lingvistică „Iorgu Iordan - Al. Rosetti”, A.R., București
Alexandru Ceaușu, Institutul de Cercetări în Inteligență Artificială, A.R., București
Lucian Chișu, Muzeul Național al Literaturii Române, București
Constantin Ciubotaru, Institutul de Matematică și Informatică, Academia de Științe a Moldovei, Chișinău
Svetlana Cojocar, Institutul de Matematică și Informatică, Academia de Științe a Moldovei, Chișinău
Alexandra Cornilescu, Facultatea de Litere, Universitatea București
Dan Cristea, Facultatea de Informatică, Universitatea „Al. I. Cuza” și Institutul de Informatică Teoretică, A.R., Iași
Nicolae Curteanu, Institutul de Informatică Teoretică, A.R., Iași
Gabriela Czibula, Universitatea „Babeș-Bolyai”, Cluj-Napoca
Cristina Florescu, Institutul de Filologie Română „Al. Philippide”, A.R., Iași
Corina Forăscu, Facultatea de Informatică, Universitatea „Al. I. Cuza”, Iași și ICIA, A.R., București
Maria Georgescu, ISSCO / TIM, ETI, Université de Genève
Gabriela Haja, Institutul de Filologie Română „Al. Philippide”, A.R., Iași
Catalina Hallett, Open University, London
Sanda Harabagiu, Department of Computer Science University of Texas at Dallas
Maria Husarciuc, Centrul de Studii Biblico-Filologice „Monumenta linguae Dacoromanorum”, Universitatea „Al. I. Cuza”, Iași
Adrian Iftene, Facultatea de Informatica, Universitatea „Al. I. Cuza” Iași
Diana Inkpen, Université d’Ottawa
Radu Ion, Institutul de Cercetări în Inteligență Artificială, A.R., București
Elena Irimia, Institutul de Cercetări în Inteligență Artificială, A.R., București
Doina Jitcă, Institutul de Informatică Teoretică, A.R., Iași
Daniel Marcu, University of Southern California
Rodica Marian, Institutul de Lingvistică și Istorie Literară „Sextil Pușcariu”, A.R., Cluj-Napoca
Dan Moldovan, University of Texas at Dallas
Vivi Năstase, EML Research, Heidelberg
Constantin Orasan, University of Wolverhampton
Oana Postolache, ISI - University of Southern California, Los Angeles
Irina Prodanoff, ILC-Pisa și Università di Pavia
Marius Pașca, Google Research
Georgiana Pușcașu, University of Wolverhampton
Vasile Rus, University of Memphis
Vitalie Scurtu, Intelligentic, Cluj-Napoca
Violeta Seretan, Departament de Linguistique, Université de Geneve
Mihai Surdeanu, University of Stanford
Valentin Tablan, University of Sheffield
Amalia Todirașcu, Université Marc Bloch, Strasbourg
Horia-Nicolai Teodorescu, Institutul de Informatică Teoretică, A.R. și Universitatea Tehnică, Iași
Dan Tufiș, Institutul de Cercetări în Inteligență Artificială, A.R., București
Ioana Vintilă-Rădulescu, Institutul de Lingvistică „Iorgu Iordan - Al. Rosetti”, A.R., București
Adriana Vlad, Facultatea de Electronică, Universitatea Politehnică București și Institutul de Cercetări în Inteligență Artificială, A.R., București

COMITETUL DE ORGANIZARE:

Vlad Alexa, FII-UAIC (vlad.alex@infoiasi.ro)

Lucian Chișu, MNLR (lucianchisu@ gmail.com)

Ioana Corăci, (ioana.coraci@yahoo.com)

Marius Corăci, Intelligentics (marius@ intelligentics.ro)

Dan Cristea, FII-UAIC și IIT-AR (dcristea@ info.uaic.ro)

Florin-Tudor Cristea, FII-UAIC (florin.cristea@ infoiasi.ro)

Lucian Gădioi, FII-UAIC (lucian.gadioi@info.uaic.ro)

Adrian Iftene, FII-UAIC (adiftene@info.uaic.ro)

Angela Ioniță, ICIA-AR (aionita@racai.ro)

Petru-Adrian Istrimschi, MediaEC-UAIC (petru.istrimschi@infoiasi.ro)

Mihai-Alex Moruz, FII-UAIC și IIT-AR (mmoruz@info.uaic.ro)

Marius Răschip, FII-UAIC (mraschip@info.uaic.ro)

Horia-Nicolai Teodorescu, IIT-AR și Universitatea Tehnică, Iași (hteodor@etc.tuiasi.ro)

Diana-Maria Trandabăț, FII-UAIC și IIT-AR (dtrandabat@ info.uaic.ro)

Dan Tufiș, ICIA-AR (tufis@racai.ro)

Eugenia Țărălungă, MNLR (eugenia12@gmail.com)

Cuprins

Cuvânt înainte	9
-----------------------------	---

Capitol I – Corpusuri vocale și prelucrarea vorbirii

Metodologie pentru constituirea și analiza unui corpus adnotat de semnale vocale - cazul SROL	13
<i>Horia-Nicolai Teodorescu</i>	
Tehnici de identificare a zonelor vocalice în secvențe rostite în limba română	23
<i>Marius Zbancioc, Horia-Nicolai Teodorescu și Monica Feraru</i>	
Aspecte metodologice de organizare a datelor și de analiză statistică a vocilor emoționale	35
<i>Horia-Nicolai Teodorescu, Ioan Păvăloi și Monica Feraru</i>	
Program de editare de contururi intonaționale în limba română bazat pe ierarhii de forme prosodice funcționale	45
<i>Vasile Apopei, Doina Jitcă și Otilia Păduraru</i>	
Corpus pentru gnatofonie: protocol, metodologie, adnotare	51
<i>Horia-Nicolai Teodorescu și Alina Untu</i>	
Dispozitiv electronic de preprocesare în regim paralel a spectrului vocal	61
<i>Mircea Hulea și Alina Untu</i>	

Capitol II – Platforme, dicționare și corpusuri pentru prelucrarea textelor

Resurse lingvistice românești în flux continuu	73
<i>Dan Cristea</i>	
Comunicarea electronică și problemele noastre ortografice – fără soluții?	81
<i>Lucian Chișu</i>	
Când migrăm la diacriticele corecte?	89
<i>Bogdan Stăncescu</i>	
Construcția automată de corpusuri multilinguale	103
<i>Tiberiu Boroș, Dan Tufiș și Alexandru Ceaușu</i>	
Parsarea comparativă a dicționarelor-tezaure românești, franceze și germane	113
<i>Neculai Curteanu, Alex Moruz și Diana Trandabăț</i>	
Realizarea unui treebank românesc	123
<i>Cenel-Augusto Perez</i>	
Monitorizarea presei în cadrul proiectului Neorom	131
<i>Ana-Maria Barbu</i>	
Emoții în cuvinte: elaborarea unei resurse multilingve	141
<i>Victoria Bobicev, Victoria Maxim, Tatiana Prodan, Natalia Burciu și Victoria Angheluș</i>	

Capitol III – Aplicații ale tehnologiilor lingvistice textuale

Sistem întrebare-răspuns antrenabil pentru limba română	153
<i>Dan Ștefănescu, Radu Ion, Alexandru Ceaușu, Dan Tufiș, Elena Irimia și Verginica Barbu-Mititelu</i>	
Un sistem de traducere automată română-franceză	165
<i>Mirabela Navlea și Amalia Todirașcu</i>	

Servicii Web interoperabile și multilinguale	175
<i>Radu Ion, Alexandru Ceaușu, Dan Ștefănescu și Dan Tufiș</i>	
Tipare hiponimice pentru limba română	185
<i>Verginica Barbu-Mititelu</i>	
Mecanismele generative ale morfologiei derivaționale	195
<i>Mircea Petic</i>	
Dezvoltarea unui parser de roluri semantice pentru limba română	203
<i>Diana Trandabăț și Dan Cristea</i>	
Supravegherea pe Internet: părerea consumatorilor despre anumite produse sau despre anumite evenimente	213
<i>Adrian Iftene, Alina-Elena Mihaila, George-Alexandru Vlad și Geta Stancu</i>	
Folosirea verbelor pentru determinarea inferențelor textuale	223
<i>Mihai Alex Moruz</i>	
Summaries / Abstracts	233
Index de autori	244

CUVÂNT ÎNAINTE

Acest volum include lucrările acceptate la a șaptea ediție a Conferinței Consorțiului de Informatizare pentru Limba Română (ConsILR), dedicată *Resurselor Lingvistice și Instrumentelor pentru Prelucrarea Limbii Române* și găzduită de Muzeul Național al Literaturii Române din București, în perioada 6-7 mai 2010. Organizatorii evenimentului au fost Facultatea de Informatică a Universității „Alexandru Ioan Cuza” din Iași, Institutul de Cercetări pentru Inteligență Artificială al Academiei Române din București, Institutul de Informatică Teoretică al Academiei Române – Filiala Iași, Muzeul Național al Literaturii Române din București și firma Intelligentics.

În contextul unor eforturi internaționale tot mai susținute pentru realizarea de resurse și instrumente lingvistice interoperabile (a se vedea în acest sens proiectele CLARIN, FlaReNet, MetaNet, MetaNet4u, Panacea, MultilingualWeb, American National Corpus, Asian LanguageGrid și multe altele) conferința ConsILR a fost dedicată resurselor lingvistice și instrumentelor pentru prelucrarea limbii române (scrise și vorbite). Articolele incluse în acest volum reflectă o parte a eforturilor cercetătorilor români de a răspunde mișcării mondiale de susținere tehnologică a limbilor naturale pentru utilizarea lor în mediile de comunicare electronică.

La Conferință au participat experți lingviști, informaticieni, matematicieni și ingineri din domeniile ingineriei limbajului, lingvisticii teoretice, inteligenței artificiale, a prelucrării semnalelor și ingineriei medicale, din București, Cluj-Napoca, Iași, Chișinău și Strasbourg. Ca și la ultimele două ediții, lucrările conferinței au putut fi urmărite online, de data aceasta însă direct în Internet, prin grija MEDIAEC, Laboratorul Multimedia al Universității Alexandru Ioan Cuza, căruia îi mulțumim și pe această cale. Mulțumim de asemenea celorlalți doi sponsori ai manifestării, UPC Business și România liberă. Programul complet al Conferinței și înregistrările audio-video pot fi consultate la adresa <http://consilr.info.uaic.ro/consilr2010/>.

Un număr de conferințe invitate au focalizat interesul publicului spre probleme legate de semantică lexicală, utilizarea diacriticelor în limba română, crearea de resurse românești de interes național (cum sunt colecția de dicționare grupate sub numele Dex-online și secțiunea românească a Wikipedia). Doar o parte din aceste conferințe au fost transmise editorilor pentru a fi incluse în prezentul volum. Partea finală a Conferinței s-a axat pe probleme de lexicografie computațională, legate de realizarea marelui Dicționar Tezaur al Limbii Române în format electronic, eDTLR.

Volumul cuprinde 22 de lucrări și este structurat în trei părți: Capitol I – *Corpusuri vocale și prelucrarea vorbirii*, Capitolul II – *Platforme, dicționare și corpusuri pentru prelucrarea textelor* și Capitolul III – *Aplicații ale tehnologiilor lingvistice textuale*. Lucrările reflectă direcțiile principale de interes actual în institutele de profil ale Academiei Române, în universități și firme de informatică din țară, dar și în câteva locuri aflate în afara României. Unele articole sunt metodologice, altele ilustrează direcții noi de cercetare și un număr mare de lucrări prezintă progrese în domeniul creării de resurse lingvistice românești. Toate lucrările incluse în volum au fost extinse și corectate de către autori după Conferință, pe baza sugestiilor referenților, întrebărilor puse în timpul prezentărilor, precum și a sugestiilor editorilor. O corecție finală a

întregului volum a fost realizată de către editori. Volumul este completat de rezumatele lucrărilor în limba engleză și de un index de autori.

Ediția a 7-a a ConsILR a marcat un progres în seria acestor evenimente, configurând din ce în ce mai clar maturitatea cercetărilor dedicate tehnologiilor lingvistice aplicate limbii române și a creării de resurse dedicate. Limba română începe să se contureze ca una din limbile semnificative în privința resurselor informatice și a tehnologiilor aplicate ei. Remarcăm totodată apariția serviciilor Web lingvistice, care schimbă radical concepția asupra localizării resurselor și a aplicațiilor. Se conturează ideea că deținerea „în casă” a resurselor și instrumentelor de prelucrare, adică în calculatorul ori serverul propriu, devine lipsită de interes. Comunicațiile de mare viteză permit acum conectarea la resurse partajate aflate oriunde pe glob. Acest lucru aduce avantajul uniformității adnotărilor, aproape obligând la adoptarea de standarde, împiedică multiplicarea de versiuni diferite și face posibilă combinarea de resurse și instrumente de prelucrare din locații fizice diferite. Este îmbucurător interesul crescând al Comisiei Europene în privința finanțării cercetărilor dedicate tehnologiilor lingvistice și a bibliotecilor digitale. În articolul său care deschide numărul 9-10 al CLARIN Newsletter (martie-iunie 2010, www.clarin.eu/newsletter/), d-l Kimmo Rossi, director adjunct al Directoratului General pentru Societatea Informațională și Media al Comisiei Europene dezvăluie abundența de apeluri de proiecte în aceste direcții ce urmează a fi lansate în toamna și iarna viitoare. Nivelul finanțărilor de proiecte curente ale Uniunii Europene în direcția Tehnologiilor Limbajului depășește 70 milioane de euro și va atinge 120 de milioane de euro în 2012. Deci numai vești bune!

Editorii mulțumesc autorilor, recenzorilor și colegilor care au contribuit la apariția prezentului volum.

Iași, București, august 2010

Editorii

CAPITOLUL 1

CORPUSURI VOCALE ȘI PRELUCRAREA VORBIRII

METODOLOGIE PENTRU CONSTITUIREA ȘI ANALIZA UNUI CORPUS ADNOTAT DE SEMNALE VOCALE – CAZUL SRoL

HORIA-NICOLAI TEODORESCU^{1,2}

¹*Academia Română și*

²*Universitatea Tehnică „Gheorghe Asachi” din Iași, Iași – România*

hteodor@etti.tuiasi.ro

Rezumat

Descriem o metodologie de constituire și analiză a unui corpus reprezentativ pentru limba română vorbită. Lucrarea are un caracter în esență programatic și metodologic, dar are și scopul de a atrage atenția că un mare număr de corpusuri folosite anterior sau în prezent în inferența lingvistică sunt deficitare metodologic după cerințele actuale și pot conduce la rezultate invalide.

1. Introducere

Complexitatea și amploarea proceselor limbilor vorbite, procese influențate de un foarte mare număr de factori, face ca analizele fonetice – fie ele de natură articulator fonetică, acustic-fonetică, perceptiv fonetică, sau fonologică – să fie adesea bazate pe seturi de date a căror reprezentativitate este incertă sau insuficient caracterizată. Față de uzanțele din domeniul foneticii și dialectologiei „clasice”, din secolele trecute, standardele științifice actuale impun câteva condiții obligatorii: (i) reprezentativitate statistică; (ii) reproductibilitate; (iii) satisfacerea condițiilor pentru analiza varianței. Conform acestor condiții, foarte puține dintre studiile din secolele trecute, care la momentul realizării au fost studii de mare profunzime – rămân azi valide științific dincolo de nivelul de observații semi-empirice. Ele își pierd valoarea de stabilire de fapte fonetice, păstrându-și doar valoare orientativă, estimativă și istorică, iar aceasta nu doar din motive de limitări tehnologice (de lipsa de înregistrări de calitate), ci mai ales din motive metodologice în (dez)acord cu cerințele actuale. În această lucrare, pornim de la condițiile precizate mai sus și derivăm un set de imperative metodologice pentru constituirea unui corpus vocal adnotat, cerințe pe care le-am aplicat la constituirea corpusului SRoL¹. Condițiile le încadrăm în două categorii: cele care privesc constituirea primară a corpusului și validarea lui (Secțiunea 2 a lucrării) și cele care privesc analiza faptelor lingvistice (Secțiunea 3).

2. Metodologie de constituire a corpusului

2.1. Criterii fundamentale

Criteriul reprezentativității statistice impune: (1). Pentru un model de limbă vorbită, un număr de vorbitori suficient de mare, suficient de reprezentativ ca proporții

¹ Corpusul „Sunetele Limbii Române –SRoL” a fost realizat la inițiativa noastră în perioada 2004-2010 de un grup care include următoarele persoane (în ordine alfabetică): Monica Feraru, Mihaela Hnatiuc, Raluca Ganea, Ramona Luca, I. Păvăloi, Laura Pistol, H.N. Teodorescu, Diana Trandabăț, Alina Untu, A. Verbuță, Oana Voroneanu, M. Zbancioc (și D. Scheianu, Univ. Pitești), cu cooperarea unui număr de 53 de persoane care au contribuit cu înregistrări.

(procentaje) pentru populația care vorbește limba respectivă (unde, prin limbă, înțelegem aici fie întreaga limbă, fie un dialect, fie doar un fenomen lingvistic din limba respectivă), conform criteriilor cunoscute în analiza statistică. (2). Pentru fiecare vorbitor, un număr suficient de repetări ale unei pronunții, pentru stabilirea reprezentativității intra-vorbitor și a eliminării unor eventuale pronunții accidental deficitare. (3). Un număr de cuvinte reprezentativ pentru limbă, în sensul că setul include o proporție semnificativă de silabe frecvente ale limbii, astfel încât să fie satisfăcute condițiile pentru analiza varianței la nivelul influenței contextului asupra pronunției vocalelor și consoanelor. (4). Un set de pronunții care să reflecte o proporție neneglijabilă a proceselor limbii, printre altele tonalitatea, încărcătura emoțională și prozodia generală, chiar dacă fenomenul fonetic analizat este punctual, astfel încât să poată fi determinate influențele acelor procese. (5). O caracterizare în detaliu a vorbitorilor, astfel încât să se respecte condiția de verificabilitate (reproductibilitate a analizei) și să se poată ulterior face analize de varianță (de determinare a influențelor diverșilor factori asupra limbii vorbite – v. mai jos).

Criteriul reproductibilității impune precizarea cu acuratețe și complet a metodologiei de culegere de date, a parametrilor instrumentarului folosit, a etapelor de prelucrare primară a înregistrărilor, a algoritmilor folosiți în prelucrare, precum și precizarea oricăror alte informații necesare pentru reproducerea ulterioară a analizei sau constituirea unui corpus echivalent. De asemenea, criteriul reproductibilității impune „publicarea” (facerea publică) a tuturor datelor, începând cu înregistrările și datele subiecților și terminând cu protocoalele utilizate și instrumentele proprii dezvoltate pentru prelucrarea și analiza datelor. Un aspect esențial al reproductibilității îl constituie precizarea tipului de vorbire, conform unui număr de criterii precum: voce educată (cultă) / needucată, voce profesională / neprofesională, monolog / convorbire, voce controlată / necontrolată, caracteristici naturale sau simulate (de ex., încărcătura emotivă), contextul socio-profesional al vorbirii (de ex., conform tipologiei de texte și comunicări precizată – ne-exhaustiv – de (Turculeț, 2002, p. 76): „monolog, dialog, povestire, interviu, știri, expunere, dispută; discuție în familie, discuție amicală, discuție particulară, discuție publică particulară (în pauze în instituție, la piață, pe stradă etc. – loc?), discuție oficială, discurs public, pledoarie juridică” etc. (citare prelucrată). Este indubitabil că, până la crearea de instrumente de adnotare automate, cu performanțe mult mai mari decât cele din prezent, realizarea unui corpus care să cuprindă toate aceste tipuri de voci, cu adnotări corespunzătoare, este foarte dificilă. Ca urmare, este important să fie prezentate clar în corpusuri constrângerile de realizare (deci, tipurile de înregistrări, între altele în raport cu clasele de mai sus, cu toate detaliile posibile). Numai în acest fel se va putea realiza, de ex., un studiu ipotetic dar util, precum „O analiză ,încrucișată’ asupra modificării triunghiului vocalelor în silabele accentuate în discursurile publice științifice față de pledoariile juridice și asupra diferențelor ce apar în funcție de sexul vorbitorului în variațiile dintre discurs științific și pledoarie”.

Criteriul completitudinii (de satisfacere a condițiilor pentru analiza varianței) impune cunoașterea în primul rând a tuturor factorilor despre vorbitor – factori familiari – precum limba maternă a mamei, locul unde a copilărit – educaționali (locul unde a urmat școala primară, nivelul maxim de educație atins, profesia – care, ultima, influențează familiaritatea cu vocabularul – gradul de educare / cultură a limbii folosite), factori sociali, medicali (foarte puține corpusuri dau informații asupra

factorilor medicali, care pot influența pronunția, dar și folosirea mai largă a limbii, de exemplu ambitusul vocal sau expresivitatea emotivă) etc. În al doilea rând, analiza varianței presupune ca în corpus să se regăsească toți factorii importanți de variabilitate a pronunției, pentru a se putea face o analiză a cauzelor ce produc o anumită pronunție² și a distinge între diversele influențe. Deci, orice corpus trebuie să includă informații complete despre trei categorii: (A) despre vorbitori; (B) despre contextul și tipul vorbirii (monolog, privat, înregistrare de laborator etc.), precum și data realizării înregistrării; (C) despre (C1) tehnica înregistrărilor, (C2) criteriile de validare și acceptare a înregistrărilor; (C4) modul de pre-prelucrare a semnalelor, (C5) modul de segmentare și (C6) adnotare, (C7) modul de extragere a caracteristicilor și (C8) de validare și eliminare a valorilor „anormale” ale caracteristicilor; (C9) modul de prelucrare statistică a datelor, inclusiv instrumentele utilizate³; (D) Analiza lingvistică a cuvintelor, propozițiilor, frazelor înregistrate. Pentru fiecare factor implicat, sunt necesare minim 5, preferabil 10 înregistrări pentru fiecare sex (fiecare cu minim trei repetări).

Criteriul eticii științifice impune păstrarea securizată a datelor personale ale vorbitorilor, informarea și acordul scris al vorbitorilor pentru a participa la teste și a face publice înregistrările, precum și validarea de către un for competent a reflectării cerințelor eticii în producerea și analiza corpusului.

Din câte știm, SRoL este până în prezent singurul corpus care satisface toate aceste criterii pentru limba română – și unul dintre puținele pe plan internațional.

2.2. *Vorbitorii și fișa vorbitorului*

Pentru fiecare vorbitor, s-a întocmit o fișă a vorbitorului, fișă care conține toate datele necesare pentru caracterizarea socio-educațională și medicală a vorbitorului. Aceste fișe sunt publice, ca și restul corpusului de voce, pe situl „Sunetele limbii române”, (SRoL). Ca urmare, oricare rezultate publicate de grupul nostru asupra limbii române vorbite, pe baza prelucrării materialului din cadrul SRoL, pot fi verificate de oricine și deci validate sau invalidate de către alte echipe de cercetare, conform cerințelor metodice actuale – spre deosebire de marea majoritate a lucrărilor curent publicate, pentru care datele sunt neverificabile (nepublice). În plus, SRoL poate fi complementat de alte grupuri de cercetare cu propriile date și utilizat în conjuncție cu alte corpusuri de voce pentru cercetări mai ample.

Fișa include date despre vârstă, sex, zonă geografică în care a copilărit și s-a format ca vorbitor, zona în care a făcut studiile primare, care se știe că fixează în mare măsură varianta dialectologică a limbii vorbite, date despre studiile universitare și locul

² Subliniem printr-un exemplu relevanța acestui deziderat. Să presupunem că se dorește analiza diferențelor dintre pronunțiile, pentru o vocală dată, în hiatus, ca vocală glisată și respectiv ca diftong. Dacă numărul de vorbitori din corpus este redus, de exemplu 6, iar ei au un fundal educațional-familial foarte diferit (dialectal, social), cu implicații în diftongizarea vocalei respective, fundal nedocumentat, în plus dacă printre ei unul are afecțiuni neurologice (nedocumentate) care reduc viteza de reacție (de ex., viteza impulsului electric pe nervi mai mică decât normal), cu implicații în producerea glisandoului vocalic, rezultatele statistice vor indica, indiferent de faptul lingvistic, real sau nu, al diftongizării glisantelor în limba dată, o „tendență de diftongizare”. Într-o asemenea analiză ipotetică, criteriul de satisfacere a condițiilor pentru analiza varianței nu este satisfăcut, iar cauzele personale ale tendinței de diftongizare sunt „văzute” ca tendință a limbii.

³ Atunci când instrumentele nu sunt publice sau comerciale, de exemplu când sunt dezvoltate de autorii corpusului, instrumentele precum și algoritmi care stau la baza lor, preferabil și codul sursă trebuie făcute publice, pentru a îndeplini condițiile de verificabilitate și reproductibilitate.

absolvirii, eventualele studii post-universitare, obiceiuri de viață care pot influența vocea (de ex., fumatul), date biometrice și patologii cunoscute. Detalii au fost prezentate mai pe larg în alte lucrări. Unul singur dintre vorbitori are patologii cunoscute care să îi afecteze sistemul fonator, respirator, auditiv, sau nervos. Mulți vorbitori au vocea „educată”, fie prin profesia didactică, fie prin obișnuința unor prezentări publice.

Studiul s-a făcut cu respectarea în totalitate a păstrării privațiunii subiecților; numele subiecților este cunoscut doar realizatorilor sitului și nu este menționat în fișa publică a vorbitorului, unde identificarea se face printr-un cod numeric. Toți vorbitorii au fost informați asupra obiectivelor generale ale cercetării și condițiilor de difuzare a datelor primare (înregistrările vocale); toți vorbitorii și-au dat consimțământul scris pentru înregistrări. Cercetarea a fost avizată sub raport etic de Consiliul Facultății ETTI, Universitatea Tehnică „Gheorghe Asachi” din Iași.

2.3. Statistica regională și socio-educatională a vorbitorilor

Setul de 53 de vorbitori include cca. 65% bărbați și cca 35% vorbitori feminini. Grupa de vârstă reprezentată este majoritar 20-35 de ani (peste 50%) și doar cațiva vorbitori sunt în grupa de vârstă 45-60 ani. Cu rare excepții, toți vorbitorii sunt educați și au copilărit în zona de NE a Moldovei, din București, în județele Iași, Vaslui, Bacău, Botoșani, Suceava. Trei vorbitori sunt din Transilvania, unul din zona Argeș, doi din Vâlcea, doi din Muntenia - București, iar unul din Maramureș. Ca urmare, SRoL permite comparații între vorbitorii (de limba cultă) din NE Moldovei cu cei din alte regiuni ale României. Profilul socio-educational al vorbitorilor este, pentru marea lor majoritate, acela al tânărului născut în anii 1980-1990 și educat până la nivel de facultate, mulți având un masterat în informatică, lingvistică, sau în domenii conexe.

2.4. Protocolul de înregistrare

Înregistrările au fost efectuate de mai mulți membri ai echipei SRoL, folosind programul GoldWave™ 5.0, cu următorii parametri: frecvență de eșantionare de 22050 Hz; rezoluție de 16 și 24 biți (utilitarul Praat™ prelucrează numai fișierele pe 16 biți), monofonie. Programul GoldWave este un program comercial, suficient de performant la nivelul anului scrierii lucrării, care asigură o înregistrare de calitate, cu parametrii doriți, precum și o prelucrare primară a datelor la nivel acustic.

2.5. Metodologia de validare a fișierelor de semnal vocal și de segmentare precisă

Toate fișierele au fost analizate manual, prin ascultare, de echipa SRoL. Cu această ocazie, au fost făcute prelucrări primare, anume au fost eliminate secțiunile zgomotoase sau pauzele prea mari; eventual, întreaga înregistrare a fost rejectată dacă nu satisfacea cerințele pentru o înregistrare de calitate, de ex. dacă prezintă „limitări” (trunchieri în amplitudine) de semnal, sau dacă avea zgomot intens, ușor perceptibil auditiv.

Două metode de *segmentare* sunt implicit sau explicit utilizate în SRoL. Prima metodă, numită aici *metoda perceptivă*, constă în ascultarea de către o persoană familiarizată atât cu fonetica generală a limbii respective, cât și cu instrumentul informatic utilizat. Metoda poate fi considerată a fi o tehnică specifică foneticii perceptive, cu limitele ei – în principal percepția dependentă de contextul lingvistic în care ascultătorul evaluator s-

a dezvoltat (limba natală, educația, cultura lingvistică). A doua metodă⁴, dezvoltată semi-empiric de autor de-a lungul anilor este subordonată în egală măsură foneticii acustice și teoriei semnalelor, deși este aplicată manual și subiectiv.

Segmentarea fișierelor de voce a fost realizată manual, spre deosebire de alte studii similare (vezi de ex. (Gendrot & Adda-Decker, 2005)), pentru a se asigura o acuratețe și un grad de certitudine cât mai mare și o compatibilitate mare între criteriile perceptive și cele utilizate în abordarea acustic-instrumentalistă. Segmentarea „de precizie” a fost realizată prin ascultare, conform limitelor percepute între foneme, cu teste repetate pentru a determina cât mai precis granițele detectate auditiv între fonemele adiacente. Atunci când ascultarea nu permitea o precizare suficient de netă în timp a granițelor, s-a analizat suplimentar forma de undă și s-a decis plasarea graniței acolo unde, vizual, apărea clar trecerea de la o formă de undă a unui sunet la forma de undă a următorului. Atunci când a fost necesar, în special pentru pauze și zone de tranziție vocală-consoană, o decizie s-a luat în grup pentru validarea segmentării.

2.6. Metodologia de adnotare

Adnotarea a fost realizată manual, simultan cu etapa de segmentare. Nivelurile de adnotare au fost cele corespunzătoare segmentelor de tip propoziție – cuvânt – silabă – fonem – segment central de fonem. Au fost distinse la același nivel cu fonemele trei tipuri de pauze: pauze între propoziții, pauze între cuvinte, pauze în interiorul cuvântului⁵. Frazele au fost adnotate de echipa SRoL folosind utilitarul Praat™, versiunea 5.1.30 (Boersma, 2002).

La nivel de fonem, segmentele au durate larg variabile, cele mai mici de ordinul câtorva milisekunde (ms), pentru pauze scurte și pentru plozive, respectiv până la ordinul zecilor de ms, pentru vocale și fricative prelungite, accentuate. În starea emoțională „tristețe”, duratele monoftongilor sunt sensibil mai mari (cu până la 50%). Variabilitatea duratelor, între foneme și la același fonem funcție de stare, de locul în cuvânt (silaba accentuată fiind tipic mai lungă) etc. face ca „populația de ferestre” de analiză să aibă dispersii mari, ca număr de ferestre per fonem.

2.7. Descriere generală a conținutului corpusului

Sumar, corpusul conține înregistrări de laborator, controlate, cu zgomot redus, de voci în general culte dar nu cultivate, fără voci profesionale din categoriile artiști, reporteri, avocați, dar cu câteva voci de persoane din învățământ, cu vorbitori fără patologie cronică sau acută, cu vârsta dominantă 20-35 de ani. Propozițiile sunt pronunțate cu ton neutru sau cu ton emotiv, iar emoțiile au fost simulate și auto-stimulate. Corpusul conține vocale susținute, izolat pronunțate, diftongi, triftongi, hiatusuri, consoane, cuvinte izolate și propoziții scurte sau de lungime medie. Specific corpusului este

⁴ Metoda constă în urmărirea vizuală a formei de undă și separarea segmentelor vocalice pe baza apariției unei deosebiri substanțiale în formele de undă și în spectrul acelor segmente, respectiv pe menținerea în cadrul aceleiași segment a zonelor în care forma de undă și spectrul își păstrează un grad mare de similitudine cu zonele anterioare. Această metodă, aplicată vizual de către expertul care realizează segmentarea, a fost dezvoltată într-o metodă obiectivă (instrument informatic) (Teodorescu, 2010 b), folosind patternuri (seturi de trăsături acustice) și distanțe definite pe spațiul acestor patternuri.

⁵ pauzele intra-cuvânt, nu neapărat între silabe, notate \$, pauzele inter-cuvinte: blanc; pauzele de scurtă durată, determinate pe semnal, dar care nu se percep, notate %.

inclusiunea pe de o parte a unor structuri gramaticale particulare, precum apozitia și subiectul dublu, iar pe de altă parte a propozițiilor ponunțate cu încărcătură emoțională. Numărul de pronunții de propoziții cu încărcătură emoțională specificată este, pentru o subclasă (descrișă mai jos) de propoziții, aproximativ egal cu numărul de pronunții cu ton neutru (fără încărcătură emoțională). În acest fel, statistica realizată este echilibrată între tonul neutru (unic în cazurile studiilor anterioare asupra vocalelor limbii române) și pronunțiile emoționale. Ca urmare, statistica prezentată are o mai mare rigoare în privința variabilității pronunțiilor posibile și deci în privința modelului complet al limbii vorbite și, în același timp, capătă o dimensiune suplimentară, a încărcăturii emoționale.

Corpusul conține un set de propoziții cu încărcătură semantică imprecis delimitată și lăsată la latitudinea interpretării vorbitorului. De exemplu, una dintre propoziții poate fi interpretată ca interogativă (în scriere marcată prin ?), interogativ-exclamativă (?!), sau ca rămasă în suspensie (...). Alte propoziții suferă interpretări de tipul simplu afirmativ (.), exclamativ (!), sau pot fi interpretate în suspensie (...). În acest fel, prin alegerea cu grijă a propozițiilor, s-a asigurat nu doar libertate de interpretare vorbitorului – și implicit un grad ridicat de naturalețe și expresivitate – dar s-a asigurat și o ponderare implicită rezonabilă între diversele prozodii și încărcături emoționale, aspect original și specific corpusului nostru. Modelul de limbă este mai realist astfel, chiar dacă analiza nu include decât un timp total relativ mic de vorbire.

Studiile viitoare vor putea corecta și crește precizia modelului de limbă vorbită, în principal prin analiza unor înregistrări de durată totală mare, echilibrate conform principiilor formulate în (Teodorescu, 2010a), anume: echilibru între pronunții cu tonalități diferite (neutru afirmativ, neutru exclamativ, neutru interogativ, neutru în suspensie) și cu activări (emoționale) diferite, pentru variate emoții. Acele studii viitoare vor trebui să implice zeci de ore de înregistrări pentru a atinge un grad mai mare de semnificație statistică față de studiul actual⁶.

2.8. Tipuri de contexte pentru vocale

Pentru a satisface criteriul completitudinii (de satisfacere a condițiilor pentru analiza varianței) este necesar ca în corpus să se afle, pentru fiecare proces studiat și pentru limbă în ansamblu, un număr suficient de cazuri care reflectă condiții diferite. De exemplu, printre condițiile care influențează pronunția vocalelor se numără contextul în care apar vocalele (avem în vedere aici, ca exemplu, determinarea unei statistici pentru triunghiul vocalelor în limbă). Este cunoscut că formantul zero (fundamentală) unei vocale este influențat semnificativ de tonalitatea propoziției, de emoția exprimată, de nivelul de accentuare al cuvântului, de contextul cuvântului și de locul unde se află plasată vocala în cuvânt, influențe reflectate de linia prozodică specifică tonalității și imprimată propoziției ca atare, prin urmare și vocalei. În același timp, este cunoscută influența „contextului imediat”, constituit de fonemul (monofonemul) anterior sau ulterior vocalei, de existența unui diftong sau de glisare, asupra formanților, în special asupra lui F_1 și F_2 . Pentru fiecare vocală, numărul de ocurențe în diversele contexte (în

⁶ Până la apariția unor instrumente a căror precizie de segmentare și adnotare automată (și interpretare) să fie echivalentă cu cea a experților umani (uneori constituiți în echipă, ca în cazul unor analize de finețe făcute de noi), asemenea studii sunt greu de realizat. Într-adevăr, activitatea la corpusul SRoL până în prezent este echivalentă cu peste 2-3 (ani × om), iar un studiu comparabil pentru zeci de ore de înregistrări ar presupune zeci de (ani × om) de activitate.

care apare vocala) semnificative statistic pentru limbă trebuie să fie suficient de mare pentru o inferență statistică. În prezent, corpusul acoperă cca. 35% din limba română, la nivelul silabic⁷ (probabilitatea de regăsire a unei silabe din limba română în corpus este 35%, altfel spus, suma probabilităților silabelor din corpus în l.r. este aproape de 0.35).

Un aspect important al contextului vocalelor este structura prozodică în care sunt pronunțate. În acest sens, distingem contexte prozodice din care vocala face parte, de tip DU (creștere F_0), F_0 constant, respectiv UD (descreștere F_0), eventual primul și ultimul context cu atributul suplimentar „creștere / descreștere rapidă”. Reprezentativitatea SRoL sub acest aspect nu a fost încă determinată.

Statisticile efectuate „încrucișat”, „în diagonală” (pentru selecții de contexte în care apar vocalele și pentru grupuri de vorbitori) se realizează în aplicațiile dezvoltate de noi și de colectivul SRoL prin selectarea acelor instanțieri de vocale care corespund criteriului respectiv. Asupra lor revenim în altă lucrare transmisă la CONSILR2010.

2.9. Tipuri de propoziții, intonații și încărcături emoționale

Statistica a fost realizată pe un set de propoziții prefixate, selectate, după cum s-a mai spus, pentru a satisface riguros criteriul de „compatibilitate multi-emoțională” (Teodorescu 2002 a). Anume, criteriul impune compatibilitatea propoziției respective cu interpretări emoționale pentru toate cele trei emoții selectate drept relevante pentru corpus (bucurie, furie, tristețe), alături de tonul neutru⁸. Alegerea stărilor emoționale a fost realizată după teste preliminare (împreună cu dr. M.S. Feraru) asupra unui set de șapte emoții; selecția s-a făcut astfel încât separarea între emoții să fie suficient de netă (ambiguitate redusă la recunoaștere, deci recognoscibilitate rezonabil de ușoară de către ascultători neutri) și în același timp, încărcarea emoțională a propoziției să fie suficient de facil de realizat pentru vorbitori. Relevanța pentru corpus a tipurilor de emoție este determinată de antagonismele bucurie – tristețe și bucurie – furie, de caracterul intens al acestor emoții, precum și de ușurința cu care pot fi exprimate suficient de distinct una de alta. Propozițiile incluse în această statistică sunt toate de natură afirmativă; faptul că nici una nu include negații este probabil limitarea principală a SRoL în prezent.

Gradul de încărcare emoțională pentru fraze a fost determinat cu ajutorul unei aplicații *online* realizată de colectiv (Pistol & Teodorescu, 2010), aplicație cu ajutorul căreia au fost colectate opiniile asupra recognoscibilității emoției în pronunția respectivă de la un număr mare de evaluatori. Conform scării folosite de (Beller, Obin, Rodet, 2008), „gradul de activare” (a emoției), pentru pronunțiile de propoziții din corpusul SRoL este între inactivare și activare mare, pe o scară cu patru trepte (inactivare, activare mică, medie și mare). Ambele tipuri de activare a emoției, activare extrovertită și introvertită sunt prezente în pronunțiile de propoziții din corpus. Toate emoțiile au fost de tip simulat („acted emotions”), dar „trăit”, în sensul că li s-a recomandat subiecților să se

⁷ Probabilitățile silabelor au fost preluate din modelele de limbă determinate pe corpusuri ample în Adriana Vlad et al., *Limba română scrisă ca sursă de informație*, Ed. Paideia, Buc., 2003

⁸ Propoziții și interpretările posibile (//): Aseară (. // ! / ...), Vine mama (. // ! / ...), Ai venit iar la mine (! / ...), Cine a făcut asta (? / ?! / ... /), Omul meu îl lucră (. // ! / ...), Oricum, îți poți câștiga locul dorit (. // ! / ...). Fraze cu subiect dublu sau apoziție: Vine ea mama!, „A trecut el așa un răstimp”, (M. Sadoveanu), O ști el careva cum să rezolve asta, Mama vine și ea mai târziu, Mama știe ea ce face, Chiar știe el ce face?

gândească la o situație reală care ar elicită emoția respectivă și deci pronunția emoțională dorită⁹.

3. Metodologia de analiză

3.1. Metodologia de prelucrare și analiză formantică (acustică)

După analiza preliminară a calității înregistrărilor, acestea au fost filtrate, astfel încât să rămână doar spectrul de frecvențe de interes în analiza formantică, anume între 70 Hz și 10 kHz. Filtrarea s-a realizat folosind facilitatea utilitarului GoldWave™ de a se preciza de către utilizator banda de filtrare și forma caracteristicii filtrului.

Pentru prelucrări spectrale (detectia formațiilor și a frecvenței fundamentale), s-a selectat un instrument informatic dintre cele mai puternice la ora actuală și în același timp gratuit, utilitarul Praat™, conceput și pus la dispoziție comunității internaționale de către P. Boersma și D. Weenink (Boersma, 2002). Analiza minuțioasă a influenței deplasării ferestrei¹⁰ asupra rezultatelor (Teodorescu, Feraru) a indicat că alegerea este convenabilă și nu produce erori la variații mici ale pasului de deplasare sau ale dimensiunii ferestrei. Utilitarul Praat furnizează valorile frecvențelor medii ale formațiilor F_0 [Hz], F_1 , F_2 , F_3 , F_4 [Hz], pentru zonele delimitate la segmentare, precum și valori instantanee ale formațiilor pentru fiecare „fereastră” de analiză.

3.2. Eliminarea erorilor

Utilitarul Praat™, deși printre cele mai performante în prezent, produce cca. 10% erori evidente, din numărul total de determinări de valori medii, la valorile primilor doi formați și aproape același procent de erori la determinările de frecvență fundamentală. Deoarece valorile „evident eronate” la care ne referim sunt adesea mult mai mari decât cele credibile, aceste erori afectează semnificativ, chiar la frecvențe de apariție de ordinul 5-10%, rezultatele statisticilor. Ca urmare, este necesară eliminarea erorilor grosiere produse de instrumentul de analiză înainte de realizarea statisticii. Dat fiind numărul mare de fișiere, nu se pune problema determinării manuale a valorilor corecte. Soluția la care am recurs este de a elimina valorile „evident” eronate precum și pe cele „aberante” (outliers). Statisticile s-au realizat doar pe valorile „valide”. Această eliminare a valorilor „anormale” poate conduce la eliminarea, pe lângă valorile eronate și a unor valori reale, dar anormale, nespecifice populației globale. Aceasta este corect din punctul de vedere al cercetării dacă suntem interesați de „procesele medii”, de caracterizarea globală și nu de elemente specifice, rare de pronunție, care pot să fie mascate de eliminările operate.

În etapa de prelucrare intermediară a datelor, au fost corectate două tipuri de erori, primul produs de instrumentul de analiză Praat™, iar al doilea este posibil a se datora

⁹ Gradul de activare al emoțiilor îl determinăm astfel: cuartila superioară (conform evaluării de către ascultatori) este „foarte expresiv”, următoarea cuartilă (50-75%) este „expresivă”, a treia cuartilă (25-50%) este puțin expresivă, iar ultima slab expresivă (ne-expresivă). Facem distincția între *expresivitățile pentru fiecare emoție în parte*, deoarece unele persoane pot activa și exprima ușor o emoție, de exemplu bucuria, dar nu și altele, de ex. „triste-țea”. O persoană cu scoruri printre primii 25% dintre vorbitori (scorurile sunt calculate astfel: 4 pct. = f. expresiv, 0 = puțin expresiv; scor personal între 0 și 16) o considerăm extrovertită, iar între 25%-75%- persoană medie.

¹⁰ Instrumentul Praat a fost setat să lucreze cu ferestre de analiză de 0,025s și cu suprapuneri ale ferestrei alungătoare (suprapuneri între ferestre succesive) de ordinul a 40%, ceea ce corespunde la deplasări de 100 ms ale ferestrei.

instrumentului sau altor factori (înregistrări deficitare, pronunții deficitare). Aplicația informatică dezvoltată permite eliminarea automată a valorilor eronate date de Praat™: (a) nedefinite pentru F_0 , acolo unde există sunet vocalic; (b) valori pentru F_0 acolo unde nu este sunet vocalic; (c) valori anormal de mici sau de mari determinate de Praat pentru F_0 , F_1 , F_2 . Anterior efectuării calculelor statistice pentru fiecare formant în parte, din fișierul cu date primare sunt eliminate valorile evident eronate. După această operație de curățare, se determină statistica primară – valoarea medie și varianța pentru fiecare formant – și se elimină valorile „aberante” (outliers), din afara intervalului $\bar{x} - 3\sigma, \bar{x} + 3\sigma$. În final, pentru valorile rămase, sunt eliminate „cozile” de 3% din populație. Statistica formantului fonemului se calculează cu valorile rămase după acest șir de corecții. Metoda este prezentată pe larg în alte lucrări.

3.3. Metoda de validare a expresiei emoționale

Validarea expresiei emoționale a fost realizată în două etape. În prima fază, cinci evaluatori din cadrul colectivului nostru au efectuat o validare a emoției reprezentate (exprimate voluntar) în pronunția respectivă a propoziției. Înregistrările care aveau un nivel de expresivitate inacceptabil de scăzut sau o exprimare confuză a emoției au fost eliminate în această etapă dintre înregistrări și înlocuite eventual cu alte înregistrări. În etapa a doua a validării, am recurs la validarea publică, pe Internet. S-a creat o aplicație (Pistol & Teodorescu, 2010) cu ajutorul căreia un utilizator oarecare de Internet, vorbitor de limba română, poate accesa și asculta oricare dintre propoziții și preciza dacă acea propoziție are încărcătură emoțională, dacă emoția este evidentă (nu produce confuzie, indecizie asupra tipului de emoție) și care este emoția respectivă.

4. Discuție și concluzii

Criteriile expuse au fost aplicate la realizarea corpusului SRoL. Acest corpus nu satisface decât parțial criteriul privind completitudinea necesară pentru analiza varianței, dar, prin aplicarea coerentă a criteriilor până acum și în continuare la dezvoltarea corpusului, premisele corectitudinii metodologice ale SRoL sunt asigurate. Nesatisfacerea criteriului completitudinii este relativă și este datorată numărului încă redus de pronunții pentru unele tipuri de ocurențe – de exemplu, pentru vocale glisante și hiatusuri. Dar în cazul SRoL limitarea nu este datorată, ca în cazul altor corpusuri, necunoașterii (ne-documentării) unor factori care pot influența pronunția – de la factori specifici vorbitorului la modul de înregistrare sau pre-procesare a fișierelor. La aceste din urmă corpusuri, eroarea este fundamentală și necorijabilă retroactiv.

Încheiem cu un exemplu de aplicație deja inițiată de noi. Să considerăm problema determinării „triunghiului vocalelor” ($F_{1(k)}^v, F_{2(k)}^v$) pentru limba română¹¹. La nivelul actual, această determinare presupune: determinarea „norilor” vocalelor, a elipselor de încredere, global pe întreaga populație (masculin și feminin), separat pe sexe, separat pentru vocale izolate, vocale în cuvinte izolate, în propoziții, influența emoțiilor asupra deplasării triunghiului formanților (deplasare cunoscută în literatură, pentru alte limbi),

¹¹ Din nefericire, l. română este una dintre puținele limbi europene pentru care încă nu există decât studii semi-empirice (Rosetti, Lăzăroiu, 1982), (Teodorescu et al. 1986), invalide la nivelul cerințelor prezente, privind triunghiul formanților. V. și (Teodorescu, Feraru, Trandabăț, 2005).

compararea triunghiului pentru vocale susținute cu cel obținut pe întreaga limbă (deci, ca medii ale mediilor pe vocale în diversele lor ipostaze), analiza efectului accentului în cuvânt, analiza efectului contextului (CVC, CV_, _VC, V.V, diftongului, glisării vocalei etc.), compararea cu alte limbi din aceeași familie, eventual analiza pe dialecte, profesii etc. La acestea se adaugă analiza modificării funcțiilor densitate de probabilitate și a parametrilor lor (varianță, asimetrie, aplatizare) pentru cazurile enunțate mai sus. Într-un sens, acesta este și (parțial) programul echipei SRoL pentru viitorul apropiat.

Mulțumiri. Autorul mulțumește tuturor colegilor care au lucrat la corpusul SRoL pentru nenumărate și îndelungi discuții. Activitatea la SRoL a fost parțial sprijinită de către Academia Română.

Referințe bibliografice

- Beller, G., Obin N., & Rodet X. (2008). Articulation Degree as a Prosodic Dimension of Expressive Speech. *Speech Prosody 2008, ISCA, Campinas, Brazil, May 6-9, 2008*, 681-684.
- Boersma, P.P.G. (2002). Praat, a system for doing phonetics by computer. *Glott International, Vol. 5 No. 9/10*, p. 341-345.
- Gendrot, C. & Adda-Decker, M. (2005). Impact of duration on F1/F2 formant values of oral vowels: an automatic analysis of large broadcast news corpora in French and German. *INTERSPEECH 2005, Sept. 4-8, Lisbon, Portugal*, 2453-2455.
- Pistol, L. & Teodorescu, H.N. (2010). A Note on Testing the Recognition of Emotional States and Tones in Speech. *Memoirs of the Romanian Academy, 2010 (under press)*.
- Rosetti, Al., Lăzăroiu, A. (1982). *Introducere în fonetică*. Ed. Științifică & Encicl., București, 1982.
- SRoL, Voiced Sounds of the Romanian Language Project, (autori: Teodorescu, H.-N., Trandabăț, D., Feraru, M., Ganea, R., Verbuță, A., Zbancioc, M., Hnatiuc, M., Voroneanu, O., Pistol, L., Untu, A., Păvăloi, I.), www.etc.tuiasi.ro/sibm/romanian_spoken_language/ro/voci_emoționale.htm.
- Teodorescu, H.-N., Buchholtzer, L., Poșa, C. (1986). *Comunicarea orală om-mașină*. Cap. 2 – Vocea și vorbire, și Cap. 4 – Sinteza vorbirii, Editura Tehnică, Seria „Tehnica la zi”, București, 1986.
- Teodorescu, H.-N., Feraru, M., Trandabăț, D. (2006). *Situl ‘Limba Română Vorbită’*. Lucrările atelierului Resurse lingvistice și instrumente pentru prelucrarea limbii române. Editori: Corina Forăscu, Dan Tufiș, Dan Cristea Iași, 3 nov. 2006, Editura Universității Al.I. Cuza Iași, 3-7.
- Teodorescu, H.-N. (2010 a). Noisy Speech Files Perceptual Mining for Speech and Noise Patterns, *PROMISE*, 29 martie 2010, Iași.
- Teodorescu, H.-N. (2010 b). AI Tools for Speech Analysis Applied to the Romanian Language. (Plenary paper), *ECC 2010*, 20-22 aprilie 2010, București.
- Turculeț, A. (2002). Tipuri de texte orale. In: Klaus Bochmann, Vasile Dumbravă (Eds.), *Limba română vorbită în Moldova istorică*, Volume 1. Leipziger Universitätsverlag, 2002, pp. 53-78.

TEHNICI DE IDENTIFICARE A ZONELOR VOCALICE ÎN SECVENȚE ROSTITE ÎN LIMBA ROMÂNĂ

M. ZBANCIOC^{1,2}, H.N. TEODORESCU^{1,2}, M. FERARU¹

¹*Institutul de Informatică Teoretică, Academia Română – Filiala Iași, România*

²*Universitatea Tehnică „Gheorghe Asachi”, Facultatea de Electronică, Telecomunicații și Tehnologia Informației, Iași – România*

{hteodor, zmarius}@etti.tuiasi.ro

Rezumat

Se prezintă un set de tehnici de segmentare utilizate pentru identificarea zonelor vocalice. Segmentarea este folosită ulterior de instrumentele de extragere a frecvenței fundamentale F_0 și a valorilor formaților F_1, \dots, F_4 . Comparăm precizia segmentării instrumentului propus cu cea a utilitarului Pratt, folosind fișiere adnotate cu mare precizie. Pentru reducerea timpului de rulare s-a optimizat calculul funcției de autocorelație, prin aplicarea unor algoritmi recurenți.

1. Introducere

Segmentarea automată a semnalelor vocale, recunoașterea automată a vorbirii, a limbii de proveniență, identificarea vorbitorului sunt domenii de cercetare cu vechime de câteva decenii, dar încă de mare actualitate. Faza de segmentare este importantă deoarece erorile acesteia afectează în mod direct performanțele extractorului de informații prozodice. Deși în literatura de specialitate se găsesc numeroase articole ce descriu diverse tehnici de segmentare automată (Rabiner & Schafer, 1978), (Calliope, 1989), (Rowde, 1991), problema segmentării nu este complet rezolvată, datorită cvasi-periodicității semnalului vocalic și a gradului mare de variabilitate a caracteristicilor fonemelor de la o limbă la alta.

În (Vidal & Marzal, 1990) se face o trecere în revistă asupra tehnicilor de segmentare insistând asupra metodelor fără constrângeri în ceea ce privește variația contururilor spectrale (SVF), metode ce folosesc o segmentare multi-nivel și o descompunere temporară pentru găsirea limitelor segmentelor. Tehnicile de segmentare combinate cu metodele de recunoaștere automată folosesc HMMs (Hidden Markov Models), parametri acustici, cum ar fi coeficienții MFCC (Mel Frequency Cepstral Coefficients), LPC (Linear Predictive Coding), tehnici de aliniere dinamică în timp (DTW engl. Dynamic Time Warping), etc. (Rabiner & Juang, 1993), (Esposito & Aversano, 2005).

(Juneja & Espy-Wilson, 2002) folosesc HMMs combinate cu SVMs (Support Vector Machines) pentru a detecta vocalele, consoanele finale, consoanele fricative și sonante, respectiv zonele de pauză, folosind 39 de parametri extrași din cepstru. (Matousek et al., 2003) folosește HMMs cu coeficienți spectrali MFCCs raportând procentaje foarte bune de 96% în acuratețea segmentării pe un corpus de date lingvistice pentru limba cehă. (Salam et al., 2009) propune o fuziune a două metode de segmentare a vorbirii, anume metode statistice bazate pe un algoritm de divergență și metode conexiuniste de învățare adaptivă MLP (Multi-Layer Perceptron). (Sarkar & Sreenivas, 2005) utilizează

o metodă bazată pe nivelul ALCR (Average Level Crossing Rate) pentru a detecta schimbările temporare semnificative în semnal. Sunt utilizate valori adaptive în funcție de SNR (Signal-to-noise ratio) și se compară performanța de segmentare automată cu fișiere segmentate fonetic manual. Colectivul nostru de cercetare a făcut comparații între fișierele de sunet sintetizate și fișierele de voce naturală (Teodorescu et al., 2009), folosind intervalele de timp specificate prin fișiere de adnotare, pentru a observa diferențele care apar între acestea la nivelul parametrilor extrași (durate, valori formanți, variație formanți). Corpus-ul „Sunetele Limbii Române SRoL” conține vocale susținute, fraze pronunțate cu diverse stări emoționale, sunete gnatosonice, adnotări, instrumente de analiză a semnalului vocalic, care sunt disponibile on-line (Teodorescu et al., 2005).

Pentru determinarea frecvenței fundamentale (F_0) există două tipuri de metode: metode indirecte și metode directe. Metodele directe de detecție a lui F_0 sunt: metoda impedanțimetrică (electroglotograma), metoda cinematografică (stroboscopică), metoda extracției F_0 pe baza formei de undă. Metodele indirecte implementate de către noi sunt: metode de analiză în domeniul timp (autocorelația, AMDF – Average Magnitude Difference Function), metode de analiză spectrală (metoda cepstrală, HPS – Harmonic Product Spectrum). Nici metodele directe nu pot fi considerate metode absolute, deoarece elementele elastice precum corzile vocale prezintă o vibrație amortizată.

Pentru validarea metodelor indirecte de extracție a frecvenței fundamentale este necesară o comparație între rezultatele obținute în cazul aplicării acestora și cele obținute în cazul utilizării metodelor directe. În acest scop am realizat adnotări de mare precizie folosind metoda bazată pe forma de undă. Lucrarea prezintă perfecționări ale instrumentului expus sumar în (Teodorescu et al., 2007), (Zbancioc, 2006). Scopul cercetării este implementarea unor metode de detecție de F_0 care să furnizeze rezultate mai bune decât utilitarul Praat™ sau alte utilitare existente.

2. Descrierea instrumentului de analiză prozodică

Instrumentul dezvoltat pentru extragerea informației prozodice conține mai multe blocuri funcționale corespunzătoare celor trei etape de procesare a semnalului vocal: preprocesarea, extragerea traseului intonațional pe baza valorilor frecvenței fundamentale, extragerea valorilor formanților superiori (Fig.1).

În etapa de preprocesare se realizează filtrări ale semnalului cu scopul de a elimina zgomotul nedorit și de limitare a benzii de frecvență în care se caută valorile formantice. De asemenea, în această etapă, un algoritm de segmentare permite eliminarea zonelor consonantice și a celor de pauză între rostiri și extragerea zonelor vocalice. Doar pentru aceste secvențe vocalice, care corespund vocalelor limbii române și consoanelor semivocalice are sens să fie realizată detecția lui F_0 și a formanților.

Acest instrument de analiză este un sistem hibrid neuro fuzzy prin blocul decizional ce determină F_0 , prin ponderarea rezultatelor furnizate de cele patru metode diferite de extragere a frecvenței fundamentale (metoda autocorelației 55%, metoda cepstrală 35%, metoda diferențelor AMDF 5% și metoda HPS 5%). Ponderile utilizate reflectă performanțele fiecărei metode, estimate ca număr de detecții eronate. Am asociat ponderi mai mici metodelor cu o probabilitate mai mare de a furniza date eronate. Detecțiile eronate (în principal prima subarmonică, respectiv primele armonici ale lui

F_0) sunt găsite prin compararea valorilor de ieșire consecutive furnizate de aceeași metodă și/sau de alte metode, realizându-se și o corecție a acestor „false” detecții de F_0 .

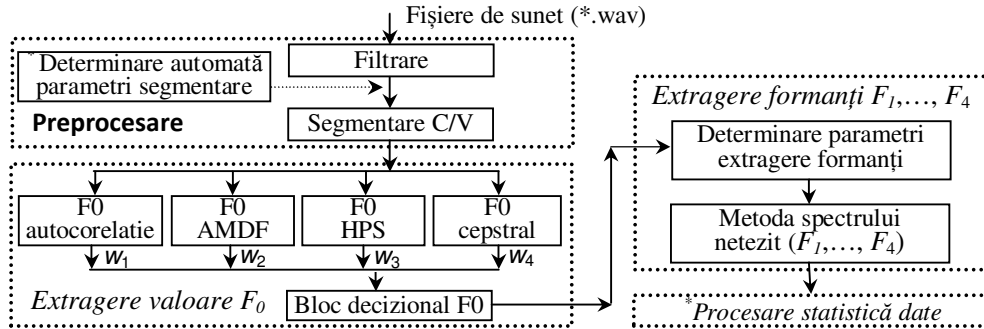


Figura 1: Schema bloc sistem hibrid neuro-fuzzy de extragere informații prozodice

Etapa de determinare a valorilor formanților superiori se bazează pe acuratețea detecției frecvenței fundamentale prin stabilirea unor intervale de căutare a fiecărui formant în funcție de F_0 . Sunt folosite tehnici fuzzy de concatenare a „spectrelor netezite” rezultate din cepstrele obținute cu diferiți parametri, dar și asocierea de coeficienți de apartenență pentru fiecare candidat găsit în benzile (intervalele fuzzy) de căutare a fiecărui formant.

Ideea utilizării unui bloc decizional hibrid este, din câte știm, originală și a permis obținerea de rezultate mai bune la o inspecție vizuală a liniei prozodice, față de alte instrumente similare de analiză a informației prozodice existente pe web: Klatt, Praat (Boersma & Weenink, 2006), Goldwave (www.goldwave.com), Wasp (www.wasp.dk), Speech Analyzer (www.sil.org/computing/sa), Winpitch (www.winpitch.com).

3. Tehnici de segmentare C/V

În această lucrare prezentăm doar prima fază de execuție a instrumentului de analiză prozodică și anume etapa de preprocesare și segmentare, care trebuie să se realizeze cu cât mai puține erori posibile, deoarece toate etapele ulterioare se bazează pe aceasta. Erorile din faza de segmentare vor afecta în mare măsură eroarea finală de detecție. Determinăm erorile de segmentare comparând segmentele vocalice identificate automat cu cele marcate manual într-un fișier adnotat de mare precizie, folosind utilitarul Praat. Metodologia de notare și marcarea a fonemelor este descrisă în cele ce urmează.

3.1 Metodologia de adnotare

Adnotarea prin ascultare este subiectivă și greu de realizat cu precizie, fiind necesară în paralel și inspecția vizuală a formei de undă și a spectrului semnalului (spectrogramei). Pe baza informațiilor legate de periodicitatea semnalului, de vârfurile spectrale și tranzițiile acestora, de componența în frecvențe înalte se poate realiza delimitarea fiecărui fonem. Chiar și utilizând toate aceste informații, stabilirea cu exactitate a acestor limite este dificilă, mai ales datorită perioadelor de tranziție dintre foneme, a perioadelor de amortizare, a zgomotelor introduse de aerul aspirat/expirat, de sunetele produse la închiderea buzelor etc.

Deoarece adnotările uzuale nu au o precizie suficientă (de ordinul ms), a fost necesară realizarea manuală a unei adnotări de mare precizie, prin metoda (practicată frecvent de al doilea autor) a comparării formei de undă în domeniul timp și domeniul frecvențelor.

Durata minimă a pauzelor intravorbire sau a secvențelor care pot fi sesizate de urechea umană este de ordin zecimi milisecunde, și de aceea pentru a se ajunge la precizia de ordin ms este necesară și analiza vizuală a semnalului. Scopul autorilor este de a determina erorile date de alte instrumente, în cazul de față Praat, și de a le compara cu cele date de extractorii de F_0 implementați de colectivul SRoL. În acest scop s-au adnotat 10 fișiere din baza de date a sit-ului SRoL (fișierele de sunet provin de la 6 vorbitori, 3 de gen feminin și 3 de gen masculin). Față de metodologia de adnotare aplicată uzual, în plus s-a ținut cont în procesul de adnotare de următoarele:

- S-au introdus notații suplimentare care să permită specificarea sunetelor specifice limbii române ('â' = 'a-', 'ă' = 'a+', 'ș' = 'sh', 'ț' = 'tz').
- S-au marcat în mod diferit zonele de pauză astfel: pauzele intervorbire (între rostiri de propoziții) cu ' ' (caracterul blank), pauzele intravorbire (silabe, cuvânt) cu '\$', pauzele care nu se aud (prezente în consoanele 'p', 't', 'c' etc.) cu '%'.
-
- S-au stabilit intervalele de demarcație ale fiecărui fonem atât prin inspecția vizuală a formei de undă (pentru a observa periodicitatea semnalelor în cazul în care acestea sunt vocalice), analiza spectrului și spectrogramei, cât și prin inspecție auditivă.
- S-a validat de mai mulți evaluatori delimitarea fonemelor și a pauzelor intravorbire.

3.2 Descriere algoritmi de segmentare C/V

Algoritmii utilizați anterior în segmentare estimau pentru fiecare fereastră de analiză (de dimensiune uzuală $N=512$ sau 1024 eșantioane) energia în domeniul timp și energia spectrală a frecvențelor joase și comparându-le cu niște valori de prag stabileau dacă acele segmente erau vocalice sau consonantice. Etapele algoritmului sunt următoarele:

- 1) Aplicarea unui filtru Butterworth de ordin 11 în banda $[70,6000]$ Hz (se păstrează informația din banda de căutare a formanților și se elimină zgomotul indus de rețea).
- 2) Parcurgerea întregului semnal și determinarea energiei maxime $E_W = \sum_{i=1}^N |s_i|$ a unei ferestre de analiză W , cu dimensiune N eșantioane. Folosind această valoare $E_{W \max}$ se parcurge din nou semnalul și se consideră că în zonele în care energia ferestrei curente este mai mică decât 20% din energia maximă nu avem zonă vocalică.
- 3) Pentru fiecare fereastră de analiză se calculează energia spectrală totală $E_{FFT}^t = \sum_{f=0}^{F_s/2} |FFT(s_W)|$ și energia din banda frecvențelor joase $[70,2500]$ Hz $E_{FFT}^B = \sum_{f=70}^{2500} |FFT(s_W)|$ (F_s este frecvența de eșantionare). Dacă $E_{FFT}^B < 0.5 \cdot E_{FFT}^t$, energia corespunzătoare frecvențelor înalte este mai mare decât 50% din energia spectrală totală, atunci se consideră că zona respectivă nu este vocalică.

Primul criteriu de segmentare a zonelor consonantice/vocalice este o metodă globală al cărui scop este acela de a elimina zonele de pauză dintre pronunții, unde este prezent doar zgomotul ambiental, precum și o serie de consoane care conțin zone de pauză (de

TEHNICI DE IDENTIFICARE A ZONELOR VOCALICE ÎN SECVENȚE ROSTITE ÎN LIMBA ROMÂNĂ

exemplu plozivele). Al doilea criteriu de segmentare C/V bazat pe o metodă locală de estimare a ponderii componentelor spectrale are rolul de a elimina consoanele care au o energie a frecvențelor înalte mai mare în raport cu frecvențele joase.

Valorile de prag utilizate în segmentare au fost determinate empiric după mai multe simulări. Algoritmii de segmentare nu funcționează bine, atunci când există zone cu energia în domeniul timp mult mai mare decât în restul pronunției. Astfel de zone apar datorită intonației (accentuarea unui cuvânt sau silabe), sau la exprimarea unei emoții (cum ar fi starea de furie, bucurie). Din acest motiv s-a testat și o variantă de segmentare în care valoarea energetică maximă este calculată pe o zonă locală restrânsă la 0.5 secunde în jurul eșantionului curent.

S-a realizat de asemenea și determinarea statistică a valorilor optime de prag folosind un algoritm de antrenare supervizat, care să determine pe baza seturilor de antrenare reguli de identificare a zonelor vocalice de cele consonantice și de zonele de pauză dintre cuvinte sau fraze. Seturile de antrenare conțineau pentru fiecare fonem: *zcr* (rata trecerilor prin zero pe secundă), *avg_e*, *std_e* - energia medie și deviația standard a energiei în domeniul timp a ferestrelor de analiză, energia spectrală în benzile B1 [70, 500]Hz, B2 [500, 1000]Hz, B3 [1000, 2000]Hz, B4 [2000, 5000]Hz. S-a preferat în locul rețelelor neuronale sau algoritmilor genetici, utilizarea arborilor de decizie, deoarece aceștia furnizează la ieșire un set de reguli, care au în premise valorile prag determinate automat pentru o clasificare cu eroare minimă. Exemplu de regulă obținută cu See5 (www.rulequest.com/see5-win.html) pentru identificarea zonelor vocalice:

```
Rule 1: (270/120, lift 1.4); aplicabilă pentru 150 pattern-uri din cele 270
IF E_MED > 0.000516
  B4 <= 0.068311
THEN class 1 (vowel) [0.555]
```

Folosirea de metode de segmentare cu valori prag nu oferă întotdeauna rezultate bune. De exemplu, aceste valori de prag nu mai pot fi folosite pentru a delimita zonele de pauză de cele în care vorbitorul rostește ceva, atunci când se vorbește încet, sau când persoana a fost amplasată prea departe de microfon (deși înregistrările s-au efectuat după un protocol care prevede o distanță optimă de la buze la microfon), sau când nivelul de zgomot ambiental este prea mare sau fonemul (vocala), de la finalul secvenței rostite are o energie scăzută și/sau o durată de atenuare mai mare. O soluție pentru această problemă este ajustarea pragurilor în funcție de SNR.

Chiar și criteriul de segmentare bazat pe energia frecvențelor înalte nu reușește să determine unele vocale aflate în hiat/diftong, dificultăți fiind întâlnite în acest caz în zonele de tranziție. Un alt caz este cel al fonemului 'a' aspirat (pronunțat în timp ce aerul este aspirat pe gură) din propoziția 'A trecut el așa un răstimp' care este încărcat în frecvențe înalte, dar își păstrează traseele formantice uzuale. O serie de foneme consonantice pot fi găsite în unele pronunții ca zone vocalice (de exemplu 'v', 'z', 'r' etc.), dar există situații în care acestea învecinate fiind cu mai multe consoane nu prezintă acea periodicitate a semnalului și sunt clasificate ca și consoane.

În consecință, este nevoie de un algoritm de segmentare care să verifice dacă semnalul este cvasi-periodic. Criteriile de segmentare utilizate anterior ar putea doar furniza informații suplimentare privind clasa de apartenență (vocală / consoană / pauză rostiri) a

semnalului din fereastra de analiză. În aceste condiții s-a preferat utilizarea funcției de autocorelație, care s-a dovedit a fi și cea mai robustă din punct de vedere a erorilor de extragere a lui F_0 , într-un algoritm adaptat pentru faza de segmentare:

- 1) Filtrarea întregului semnal folosind un filtru de mediere de ordin $N=31$;

$$s_{filt}[k] = \sum_{i=1}^{31} s[k+i]/31, \quad k = \overline{1, len_s - 31},$$

unde len_s reprezintă lungimea semnalului de intrare. Alegerea ordinului filtrului este justificată de rezultatele simulărilor, prezentate în secțiunile următoare.

for $i = 1$ to $n_iteratii$

- 2) Calcularea vectorului de valori $scorr$ prin aplicarea funcției de corelație pentru fereastra curentă de analiză.

- 3) Căutarea maximumului din vector și calcularea valorii $vF0 = Fs / poz_{max}$. Dacă această valoare este situată între $poz_{start} = Fs / 500$ și $poz_{end} = Fs / 80$, atunci „considerăm” acel segment ca fiind vocalic, altfel $vF0 = 0$.

end_for

- 4) Se determină vectorul boolean $\{z(vF0) | z : [0,500] \rightarrow \{0,1\}\}$ al variabilității semnalului $vF0$, considerând că între două valori consecutive există variabilitate, dacă între acestea avem o tranziție (variație) mai mare $\pm 5\%$.

$$z[k] = \begin{cases} 1 & vF0[k] > 1.05 \cdot vF0[k+1] \text{ \& } vF0[k] < 0.95 \cdot vF0[k+1] \\ 0 & \text{altfel} \end{cases}, \quad k = \overline{1, len_s - 1}.$$

- 5) Dacă în vectorul z avem variabilitate mai mare de $p\%$ atunci considerăm zona respectivă ca nefiind nevocalică. Se obține în final semnalul \hat{s}_{segm}

$$\hat{s}_{segm}[k] = \begin{cases} s[k] & , \hat{z}[k] = \sum_{i=-N/2}^{N/2} z[k+i] < N \cdot p\% \\ 0 & \text{altfel} \end{cases}.$$

În ultima etapă a algoritmului am considerat pragul de variabilitate de $p=5\%$: dacă din $N=500$ de valori avem mai puțin de 20 care variază cu $\pm 5\%$ față de valorile vecine, considerăm zona respectivă ca fiind vocalică. Pentru vectorul $vF0$ se pot stabili și reguli statistice care să impună ca deviația standard pentru un număr consecutiv de valori să nu depășească o valoare prag dată. Variabila $n_iterații$ se deduce în funcție de dimensiunea ferestrei de analiză, w și de pasul de deplasare al ferestrei, $n_iteratii = \lfloor (len_s - w) / step \rfloor$. Pentru o segmentare cu o rezoluție maximă a semnalului de ieșire, pasul de deplasare se alege $step = 1$.

În etapa a treia se consideră că o valoare de pe poziția k din vectorul $scorr$ este maxim local, dacă este mai mare decât valoarea din dreapta $scorr[k+1]$ și din stânga acesteia $scorr[k-1]$. Căutarea maximumului corespunzător lui T_0 se face doar în lista maximelor locale astfel găsite. Se evită astfel selectarea ca maxim a primei valori din vector $scorr[0] = \sum_{i=1}^N s^2[i]$, care corespunde energiei semnalului din fereastra curentă de analiză.

Un dezavantaj al metodei propuse îl constituie numărul mare de operații care trebuie efectuate. Sunt necesare cinci parcurgeri ale unor vectori comparabili ca dimensiune cu semnalul de intrare, s_{filt} obținut după filtrare folosind o fereastră de $N=31$ eşantioane,

TEHNICI DE IDENTIFICARE A ZONELOR VOCALICE ÎN SECVENȚE ROSTITE ÎN LIMBA ROMÂNĂ

s_{corr} rezultat după calcul funcției de autocorelație cu $N=1024$ eșantioane, z vector variabilitate pentru fiecare două eșantioane consecutive și \hat{z} vectorul variabilității cumulate pentru $N=500$ (folosit la generarea semnalului de ieșire \hat{s}_{segm}).

Pentru a obține timpi de rulare mai mici am optimizat algoritmul de segmentare prin folosirea unor funcții recursive de calcul a vectorilor s_{filt} , s_{corr} , \hat{z} , esențială fiind scăderea timpului de calcul la nivelul funcției de autocorelație.

3.3 Optimizarea algoritmului de segmentare prin funcții recursive

Funcția de autocorelație necesită un număr mare de operații, timpi mari de calcul, fiind un algoritm de complexitate $O(N^2)$, motiv pentru care s-a evitat anterior utilizarea ei în faza de segmentare pentru a studia periodicitatea semnalelor. Din această cauză, în cazul extractorului de F_0 s-a preferat folosirea unui pas mai mare de deplasare a ferestrei $step = N/16$, unde $N=1024$ reprezintă dimensiunea implicită a ferestrei de analiză.

Pentru rezoluția maximă ($step=1$) s-au obținut pentru un fișier de 10-20 secunde timpi de rulare de ordin 5-10 minute. Prezentăm pseudo-codul algoritmului de calcul a vectorului de autocorelație, pentru fiecare fereastră de analiză, W , selectată la parcurgerea semnalului de analizat, s cu pas de deplasare al ferestrei $step$.

```
for j=1: step : niter*step
//Calculează  $C_{XX}^{W_j}$  funcția de autocorelație pentru fereastra curentă  $W_j$ 
    for k=0:N
        for i=1:P
             $C_{XX}[k] = C_{XX}[k] + x[i] \cdot x[i+k]$ 
```

În limbajul MATLAB de exemplu, funcția de autocorelație folosește doar valori din fereastra curentă (de dimensiune N), motiv pentru care variabila i variază între $1: N-k$. În acest caz formula de calcul pentru funcția de autocorelație devine:

$$C_{XX}[k] = \sum_{i=1}^{N-k} x[i] \cdot x[i+k], \quad k = \overline{0, N}$$

Pentru optimizare am folosit o relație de recurență în calculul funcției de autocorelație a unei ferestre W_2 pornind de la șirul de valori al ferestrei anterioare $C_{XX}^{W_1}$.

$$\begin{aligned} C_{XX}^{W_1}[0] &= x_1^2 + x_2^2 + \dots + x_N^2 & C_{XX}^{W_1}[k] &= x_1 \cdot x_{k+1} + x_2 \cdot x_{k+2} + \dots + x_{N-k} \cdot x_N \\ C_{XX}^{W_2}[0] &= x_2^2 + \dots + x_N^2 + x_{N+1}^2 & C_{XX}^{W_2}[k] &= x_2 \cdot x_{k+2} + \dots + x_{N-k} \cdot x_N + x_{N-k+1} \cdot x_{N+1} \\ C_{XX}^{W_2}[0] &= C_{XX}^{W_1}[0] - x_1^2 + x_{N+1}^2 & C_{XX}^{W_2}[k] &= C_{XX}^{W_1}[k] - x_1 \cdot x_{k+1} + x_{N-k+1} \cdot x_{N+1} \end{aligned}$$

Relația de recurență de mai sus este calculată pentru cazul în care între cele două ferestre avem un pas de deplasare minimal ($step=1$), fapt ce conduce la o rezoluție maximă a semnalului de prelucrat. Procesul computațional ar necesita în loc de $N \cdot (N+1)/2$ operații de adunare și înmulțire un număr de doar $2 \cdot N$ operații și, teoretic, timpul de calcul ar trebui să scadă de $N/4$ ori. Astfel pentru parcurgerea întregului semnal cu o fereastră de analiză de 1024 eșantioane și calculul funcției de autocorelație am avea teoretic un timp de calcul de aproximativ 250 ori mai mic.

Crescând pasul de deplasare, crește și numărul de operații de efectuat, astfel încât timpii finali de calcul nu vor fi mai mici decât în cazul în care $step=1$. Aceasta se observă din relația de recurență rezultată:

$$C_{XX}^{W2}[k] = C_{XX}^{W1}[k] - \sum_{j=1}^{Step} x_j \cdot x_{k+j} + \dots + \sum_{j=1}^{Step} x_{N-k+2-j} \cdot x_{N+2-j}$$

Pentru filtrarea semnalului s-a preferat folosirea filtrului de mediere neponderat în locul unui filtru digital FIR Butterworth, Bessel, Chebyshev pentru a putea aplica o relație de recurență. Pentru filtru de mediere ponderat:

$$s_{filtr}[k] = \left(\sum_{i=1}^N s[k+i] \cdot a[i] \right) / \sum_{i=1}^N a[i],$$

o relație de recurență devine posibilă doar dacă coeficienții $a[i]$ sunt egali între ei ($a[0] = \dots = a[i] = \dots = a[N]$ obținându-se un filtru de mediere neponderat):

$$s_{filtr}[k+1] = s_{filtr}[k] - s[k]/N + s[k+N]/N, \quad N=31.$$

În aceeași manieră s-au folosit funcții recursive la calculul vectorului de variație \hat{z} .

$$\hat{z}[k+1] = \hat{z}[k] - z[k] + z[k+N], \quad N=500$$

Folosirea funcțiilor recurente de calcul a condus la obținerea unor timpi de rulare mai mici de câteva zeci de ori, făcând posibilă aplicarea algoritmului de segmentare propus.

4. Simulări și rezultate

În urma implementării practice a metodei de segmentare propuse, pentru obținerea unor erori minime în segmentare s-au formulat următoarele întrebări:

- unde este util să se facă filtrarea: înainte sau după calculul vectorului de autocorelație?
- care este ordinul optim pentru filtrul de mediere?
- care este valoarea de prag $p\%$ care trebuie folosită pentru obținerea lui \hat{s}_{segm} ?

Pentru a răspunde la primele două întrebări s-au folosit semnale armonice de test

$$s = A_1 \cdot \sin(2\pi \cdot f_0 \cdot t) + A_1 \cdot \sin(2\pi \cdot f_1 \cdot t) + A_3 \cdot \sin(2\pi \cdot f_2 \cdot t)$$

în care valorile frecvențelor f_0, f_1, f_2 se aleg apropiate de valorile frecvenței fundamentale F_0 și respectiv ale frecvențelor primilor doi formanți, F_1 și F_2 . În aceste condiții, valoarea de maxim care trebuie detectată și salvată în vectorul vFO ar trebui să fie cât mai apropiată de f_0 .

S-au testat două variante de filtrare, una în care se realizează filtrarea înainte (independent) de vectorul funcției de autocorelație (notată în figurile 2 și 3 cu (v2)) și una în care se realizează filtrarea în final după calcul vectorului $scorr$ (notată cu (v1)). Ne interesează acest studiu, deoarece dacă în urma simulărilor s-ar obține erori de detecție a frecvenței fundamentale semnificativ mai mici pentru (v1), decât pentru (v2) atunci nu s-ar mai putea optimiza algoritmi de segmentare folosind relații recurente. În figura 3 este reprezentată eroarea medie de detecție a lui F_0 , pentru mai multe studii de caz $\{F_0, F_1, F_2\} = \{100, 350, 900\}, \{200, 550, 1100\}, \{70, 350, 900\}, \{100, 400, 800\}, \dots$, variind ordinul filtrului în mulțimea $\{1, 5, 11, 15, 21, 25, 31, 35, 41, 45, 51\}$.

TEHNICI DE IDENTIFICARE A ZONELOR VOCALICE ÎN SECVENȚE ROSTITE ÎN LIMBA ROMÂNĂ

Din figura 2 se observă că eroarea de detecție a frecvenței f_0 este minimă pentru un filtru de mediere de ordin 31 și că algoritmul de segmentare funcționează mai bine atunci când funcția de autocorelație este calculată pentru o fereastră de 1024 eșantioane, în varianta v2, după ce în prealabil s-a realizat și filtrarea semnalului.

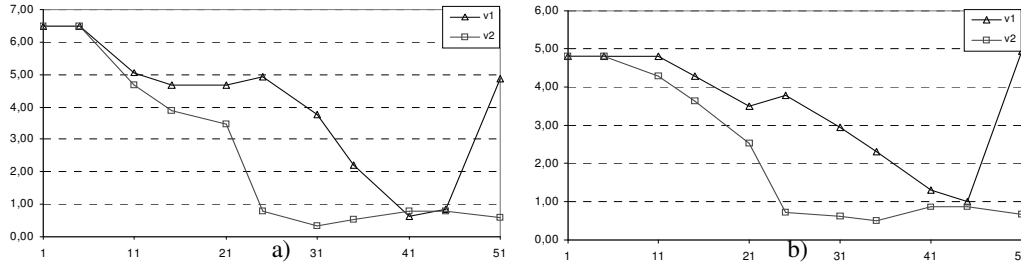


Figura 2: Eroare de detecție a lui F_0 în condițiile în care (v1) filtrarea se face după calculul funcției de autocorelație, sau (v2) filtrarea se face înainte de calculul vectorului *scorr*
 a) fereastră de analiză de 512 eșantioane, respectiv b) 1024 eșantioane

Funcționarea optimă a metodei de segmentare cu un filtru de mediere de ordin $N=31$, poate fi explicată prin comportamentul de filtru trece jos FTJ, cu valoarea riplului de aproximativ 520Hz, în condițiile în care banda de lucru este [70-500]Hz.

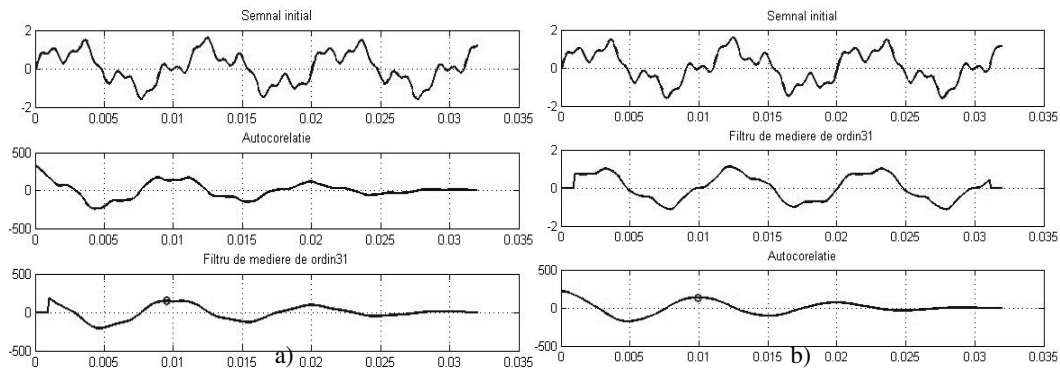


Figura 3: Semnal de test ($s = \sin(2 \cdot \pi \cdot t \cdot 100) + 0.5 \cdot \sin(2 \cdot \pi \cdot t \cdot 350) + 0.2 \cdot \sin(2 \cdot \pi \cdot t \cdot 900)$)
 a) varianta v1 filtrarea se face după calculul funcției de autocorelație, $f_0=105.3$
 b) varianta v2 filtrarea se face înainte, $f_0=100.6$

Un extractor de frecvență fundamentală furnizează următoarele erori: (i) E1: determină F_0 acolo unde nu este sunet vocalic (fals pozitiv); (ii) E2: nu determină F_0 acolo unde este sunet vocalic (fals negativ); (iii) E3: determină eronat F_0 ca valoare. (iv) E4: erori de decizie asupra valorii finale a lui F_0 , atunci când avem mai multe metode de detecție.

Erorile de segmentare automată sunt specificate de E1 și de E2. La acestea pot contribui și erorile de adnotare (segmentare manuală) în cazul în care aceasta nu este realizată cu precizie. Erorile extractorilor de F_0 sunt date de valoarea lui E3, iar la eroarea finală se poate adăuga și eroarea blocului decizional a sistemului hibrid neuro-fuzzy E4.

Semnificative ca valori sunt primele două surse de erori, de care este responsabil blocul de segmentare, E3 fiind eliminate aproape în întregime de algoritmi de corecție, iar valorile lui E4 sunt aproape neglijabile și cuantizabile doar prin inspecție vizuală.

În figurile 3 și 4 sunt prezentate rezultatele segmentării/detekției de F_0 cu instrumentele proprii și cu PraatTM pentru aceeași pronunție „O ști el careva cum să rezolve asta”) din fișierul de sunet *accent_cuv_urm_v2.wav*. Dacă algoritmi anteriori de segmentare nu

reuşeau să separe toate zonele vocalice, având dificultăţi cu acele foneme care aveau valoarea amplitudinii/energiei mai mică, cu algoritmul nou propus se segmentează foarte bine, chiar şi zone pe care instrumentul Praat nu le detectează. Astfel fonemele 'l' şi 'e' din cuvântul 'rezolve', vocala 'a' finală au traseul intonaţional detectat fără discontinuităţi. Nu avem nici false detecţii în cazul fonemelor 'ş' şi 'c'. Algoritmul nou propus are câteva zone înguste de ordin ms, situate în zonele de pauză, în care găseşte izolat valori periodice. Pentru a le elimina se pot adăuga restricţii, ca zonele considerate ca fiind foneme vocalice să aibă o durată de minim 5ms, sau restricţii privitoare la energia semnalului care în zonele de pauză sunt semnificativ mai mici decât în segmentele rostite, dacă înregistrarea nu a fost realizată într-un mediu zgomotos.

Erorile de segmentare prezentate în tabelul 1 au fost realizate considerând ca zone vocalice segmentele extrase din fişierele de adnotare corespunzătoare vocalelor, diftongilor, consoanelor sonante 'l', 'm', 'n', 'r'. Nu sunt luate în calcul unele foneme care uneori se comportă atât vocalic (forma lor de undă este periodică), cât şi consonantic. Este cazul fonemului 'v' care, în 'careva' este vocalic, dar în 'rezolve' este consonantic. Zonele de tranziţie dintre foneme vocalice, uneori sunt şi ele nevocalice (de exemplu pentru pronumele 'el' pronunţat ca regionalism 'iel', s-au găsit astfel de tranziţii între cele două vocale). Din zonele detectate de program ca fiind vocalice s-au eliminat şi la stanga şi la dreapta o jumătate din fereastră de analiză (N/2 eşantioane), corespunzătoare în general zonelor de tranziţie dintre consoane şi vocale.

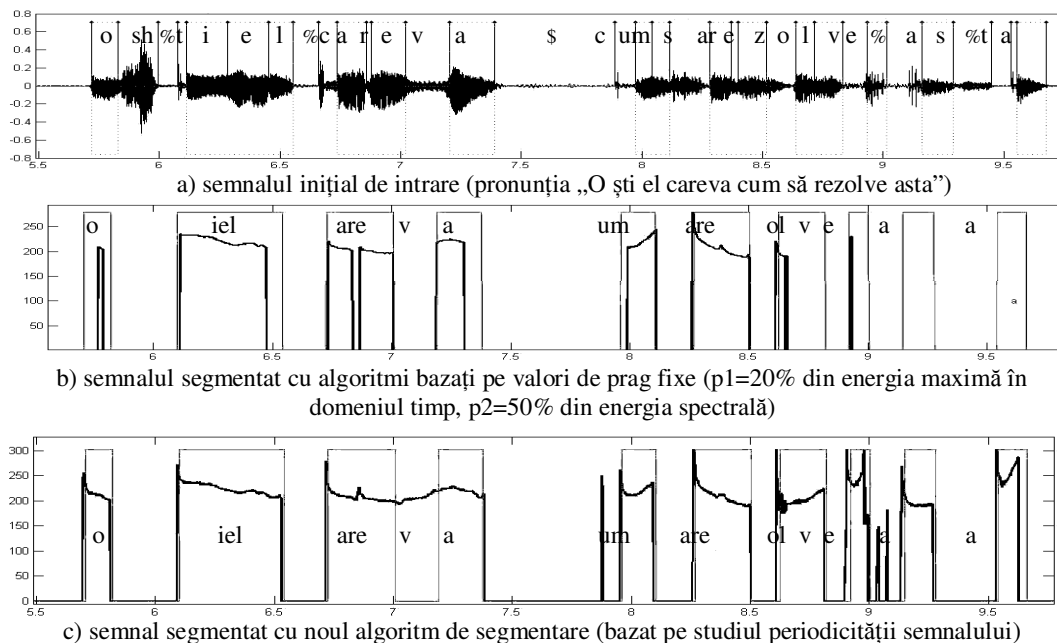


Figura 4. Rezultate segmentare și detecție F0 (instrumente proprii)

TEHNICI DE IDENTIFICARE A ZONELOR VOCALICE ÎN SECVENȚE ROSTITE ÎN LIMBA ROMÂNĂ

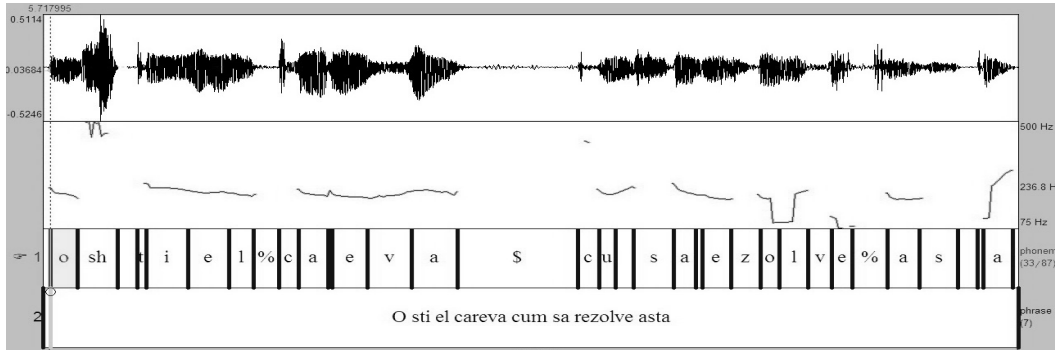


Figura 5. Rezultate segmentare și detecție F0 (soft Praat™)

Tabel 1: Erori extractor de frecvență fundamentală

Fișier de intrare	F0 - Metoda autocorelației			F0 – Praat™		
	E1	E2	E3	E1	E2	E3
accent_cuv_urm_v1.wav	0.133	0.064	0.032	0.192	0.011	0.057
accent_cuv_urm_v2.wav	0.084	0.096	0.058	0.153	0.044	0.049
accent_cuv_urm_v12.wav	0.128	0.100	0.065	0.184	0.018	0.763

5. Concluzii. Direcții viitoare

Algoritmul de segmentare propus este mai puțin influențat de amplitudinea și fluctuațiile în amplitudine ale semnalului analizat și prin urmare poate fi adaptat și pentru înregistrări cu un nivel de zgomot mai mare. Este necesară introducerea unor criterii noi de segmentare pentru detecțiile izolate de F_0 din zonele de pauză, respectiv definirea unor valori de prag flexibile, ajustate automat în funcție de SNR (va trebui estimată amplitudinea zgomotului și cea a semnalului util). Se va încerca utilizarea altor parametri care să fie mai puțin influențați de nivelul energetic, cum ar fi rata trecerilor prin zero, pentru semnalul din care s-au eliminat întâi componentele de joasă frecvență.

Conform simulărilor, algoritmi proprii oferă erori mai puține de tipul E1 (F_0 în zone nemarcate ca sunet vocalic) și de tipul E3 (se determină eronat F_0 ca valoare – cu fluctuații). În lipsa unor metode de corecție, soft-ul Praat™ produce „falsele” detecții, ceea ce reprezintă un inconvenient în eroarea globală a sistemului, în comparație cu instrumentele proprii care elimină în proporție foarte mare aceste erori.

Mulțumiri. Cercetarea a fost realizată cu sprijinul Academiei Române, în cadrul temei interne a Institutului de Informatică Teoretică din Iași. Autorii mulțumesc celorlalți co-autori ai sitului Sunetele Limbii Române, precum și referenților anonimi pentru sprijinul și observațiile pertinente.

Contribuția autorilor: Primul autor a implementat instrumentele de analiză în mediile de programare MATLAB și C++; a elaborat metoda de extragere formanți și optimizarea metodelor de segmentare și extragere a F_0 și a formanților. Al doilea autor a inițiat tema, a coordonat activitatea de cercetare, a elaborat conceptul general de sistem de decizie și metode de extragere F_0 conform cu fig.1, a precizat metoda de adnotare și segmentare manuală cu elemente de noutate privind metoda proprie de combinare a informației temporare cu cea spectrală. Al treilea autor a realizat înregistrările, a efectuat manual adnotările de mare precizie și a identificat problemele de segmentare. Toți autorii au contribuit la analiza rezultatelor și identificarea soluțiilor de îmbunătățire și validarea rezultatelor acestora.

Referințe bibliografice

- Boersma, P., Weenink, D., Institute of Phonetic Science, University of Amsterdam, Praat: doing phonetics by computer, *www.praat.org*.
- Calliope (1989). *La parole et son traitement automatique*, ISBN 2-225-81516-X, Masson, France.
- Esposito, A. & Aversano, G. (2005). Text independent Methods for speech Segmentation, *Lecture Notes in Computer Science*, ISBN 978-3-540-27441-4, 3445, 261-290 (<http://www.springerlink.com/content/81fpb3brpq367j7gf>).
- Juneja, A. & Espy-Wilson, C. (2002). Segmentation of Continuous speech using acoustic-phonetic parameters and statistical learning, *Proceedings International Conference on Neural Information Processing*, (<http://www.ece.umd.edu/~juneja/paper1910.PDF>), Universitatea din Maryland, Singapore, SUA.
- Matousek, J., Tihelka, D., Psutka, J. (2003). Automatic Segmentation for Czech concatenative speech synthesis using statistical approach with boundary-specific correction, *Processing of Eurospeech 2003*, Geneva, 301-304.
- Rabiner, L.R., Juang B.H. (1993). *Fundamentals of Speech Recognition* Englewood Cliffs, N.J.
- Rabiner, L.R. Schafer R. W. (1978). *Digital Processing of Speech Signal*, Prentice-Hall, Inc. Englewood Clifford, 11-65.
- Rowden, C. (1991). *Speech Processing*, McGraw - Hill Book Company, Chapter 2, 35-74.
- Salam, M.S., Mohamad, D., Salleh, S.H. (2009). Improved Statistical Speech Segmentation Using Connectionist Approach, *J. of Computer Science*, ISSN 1549-3636, 5 (4): 275-282.
- Sarkar, A. & Sreenivas, T.V. (2005). Automatic speech segmentation using average level crossing rate information, *IEEE Int. Conf. on Acoustics, Speech, and Signal Processing: ICASSP*, Philadelphia, SUA, 1, 397-400.
- Teodorescu, H.N., Feraru, M., Pistol L., Zbancioc, M. și alții. (2005). SRoL – Proiectul Sunetele Limbii Române (Voiced Sounds of Romanian Language Project), 2005. [http://iit.iit.tuiasi.ro/romanain_spoken_language/index.htm]
- Teodorescu, H.N., Feraru M., Zbancioc M.D. (2009) Assessing the Quality of Voice Synthesizers, In Burileanu C., Teodorescu H.N. (Eds.), *Advances in Spoken Language Technology*, The Publishing House of the Romanian Academy, Bucharest, România, ISBN 978-973-27-1808-7, 53-66.
- Teodorescu, H.N., Trandabat, D., Feraru, M., Zbancioc, M., Luca, R. (2007). A Corpus of the Sounds in the Romanian Spoken Language for Language-Related Education, In Carlos Perinan Pasqual (Eds.), *Revisiting Language Learning Resources*, Cambridge Scholars Publishing (CSP), UK, ISBN 1-84718-156-2, 6, 73-89.
- Vidal, E. & Marzal, A. (1990). A review and new approaches for automatic segmentation of speech signals, *Signal Processing V: Theories and Applications*, Torres, L., Masgrau, E., Lagunas, M.A. (eds.), Elsevier Science Publishers B.V. – Universitatea Politehnica din Valencia, Spania.
- Zbancioc, M. (2006). Tools for the Archive of the Romanian Language Sounds Project, 4th *European Conf. on Intelligent Systems and Technologies*, ECIT'2006, Iași, Romania.

ASPECTE METODOLOGICE DE ORGANIZARE A DATELOR ȘI DE ANALIZĂ STATISTICĂ A VOCILOR EMOȚIONALE

HORIA-NICOLAI TEODORESCU^{1, 2}, IOAN PĂVĂLOI¹, MONICA FERARU¹

¹*Institutul de Informatică Teoretică al Academiei Române,*

Filiala Iași a Academiei Române

²*Universitatea Tehnică Gheorghe Asachi din Iași*

{hteodor}@etc.tuiasi.ro

Rezumat

Prezentăm o metodologie și un program pentru analiza statistică a vocilor emoționale. Prezentăm în principal unele rezultate preliminare privind modul de organizare a datelor în cadrul unei aplicații de analiză a adnotărilor pentru fișierele de semnal vocal. Programul permite analiza statistică a caracteristicilor formantice ale vocalelor și semivocalelor pe subclase de fișiere selectate conform unor caracteristici date de utilizator. Aplicația permite un grad de rafinare a analizei emotivității în voce la nivelul sunetelor vocalice.

1. Introducere

Analiza expresivității emotive (Teodorescu & Feraru, 2007), (Teodorescu & Feraru, 2009) stabilirea unor metodologii teoretice și a unor algoritmi pentru identificarea stărilor emoționale (Ververidis & Kotropoulos, 2006), (Nakatsu et al., 1999), (Mcgilloway et al., 2000) au numeroase aplicații practice (Scherer, 2000), (Kienast & Sendlmeier, 2008) și științifice în domeniul informaticii, lingvisticii, psihologiei etc. Aceste analize sunt însă laborioase și necesită resurse considerabile, atunci când privesc corpusuri mari de înregistrări de vorbire. În această lucrare prezentăm realizarea unui program care efectuează o statistică generală a formanților, cu aplicație la unele vocalele („a”, „e”, „i”, „u”, „ă”) din limba română cu selecția fișierelor de voce în funcție de starea emoțională.

Scopul cercetării raportate a fost stabilirea metodologiei de organizare și de procesare a datelor în vederea prelucrării statistice. Aceasta permite flexibilitate în accesarea informațiilor, validarea datelor existente, rapiditate în prelucrarea datelor, precum și dezvoltarea de noi aplicații în domeniul procesării limbajului vorbit. Tratarea unitară a volumului mare de date ne permite o abordare eficientă în vederea realizării unor prelucrări statistice complexe. Lucrarea este o continuare a cercetărilor autorilor privitor la analiza și caracterizarea emotivității în voce. Lucrarea este structurată după cum urmează: secțiunea a doua și a treia sunt dedicate metodologiei de organizare și de procesare a datelor în aplicația dezvoltată de noi (și numită WinCollection); în secțiunea a patra sunt prezentate rezultatele preliminare obținute, iar în ultima secțiune sunt descrise concluziile și direcțiile viitoare.

2. Metodologia de organizare a datelor

2.1 Datele primare

Propozițiile înregistrate sunt: „Vine mama”, „Cine a făcut asta”, „Ai venit iar la mine” și „Aseară”. Stările emoționale supuse analizei sunt: bucurie, tristețe, furie și tonul neutru. Exemplificăm analiza pentru un număr de cinci vorbitori (2 feminini și 3 masculini) pentru care avem 63 de înregistrări cu o medie de trei rostiri pe fiecare înregistrare. Există trei înregistrări masculine (bucurie, tristețe și neutru) care au un grad de exprimare a emoției redus comparativ cu ceilalți vorbitori. Toate înregistrările fac parte din corpusul adnotat și documentat SRoL (Teodorescu et al., 2005) și au fost realizate cu o frecvență de eșantionare de 22kHz, cu o rezoluție de 16 biți.

Aplicația este reprezentată de un program realizat în Visual C++. Datele de intrare sunt conținute în patru fișiere: un fișier .wav (înregistrarea propriu-zisă) și trei fișiere .txt (un fișier de adnotare TextGrid și două fișiere cu valori ale frecvenței fundamentale, F0 și ale formanților F1-F3). Adnotarea fișierelor s-a realizat cu utilitarul Praat (Boersma & Weenink, 2006), la nivel de propoziție, cuvânt, silabă și fonem. În procesul adnotării am definit următoarele tipuri de pauze: i) pauzele intravorbire (în rostirea de silabă, cuvânt), care au fost marcate cu simbolul „\$”; ii) pauzele intervorbire (între rostiri de propoziții) care au fost codate prin blank; iii) pauzele care nu sunt percepute și care sunt detectate la nivelul unei analize mai fine; acestea au fost notate cu simbolul „%”.

Atât în fișierele Praat cât și în fișierele de intrare ale programului am codat: â cu a-, ă cu a+, ș cu sh și ț cu tz.

Fișierele text cu valori instantanee ale formanților au fost obținute cu ferestre alunecătoare a căror dimensiune este de 0.025s, iar pasul ferestrei este de 0.01s. Fișierele rulate de program au o durată de maxim 10s, ceea ce corespunde în medie la trei rostiri a propoziției înregistrate. Programul a rulat doar peste fișierele date de utilitarul Praat, dar poate prelucra și alte tipuri de fișiere.

În figura 1, exemplificăm rezultatul rostirii pentru vorbitorul 11861, pentru care avem testate 3 stări emoționale: bucurie, furie și ton neutru. Frazele ale căror pronunții sunt supuse analizei sunt: „Ai venit iar la mine”, „Cine a făcut asta”, „Vine mama” și „Aseară”, iar fonemele analizate sunt „a”, „e”, „i”, „e”, „ă”.

```

speaker : 11861
  Emotion : bucurie
    Phrase : ai venit iar la mine
              3 rostiri
    Phrase : aseara
              3 rostiri
    Phrase : cine a facut asta
              3 rostiri
    Phrase : vine mama
              3 rostiri
  Emotion : furie
    Phrase : ai venit iar la mine
              3 rostiri
    Phrase : aseara
              3 rostiri
    Phrase : cine a facut asta
              3 rostiri
    Phrase : vine mama
              3 rostiri
    
```

Figura 1: Screenshot referitor la sinteza analizei pentru vorbitorul 11861

2.2 Preprocesare de fișiere

Valorile lui F0 și F1-F4 au fost calculate automat folosind utilitarul Praat și salvate în fișierele corespunzătoare vorbitorului, emoției și propoziției respective. Există unele situații în care F0 apare ca fiind nedefinit. Aceste situații sunt următoarele:

- între rostiri - nu există F0 și nici F1-F4. În tabelul 1 exemplificăm un segment de valori ale lui F0 care este nedefinit și altul definit pentru propoziția „Ai venit iar la mine”, diftong „ai”, prima rostire. Primele valori nedefinite sunt cele care corespund valorilor inexistente din intervalul dinaintea rostirii. Programul ignoră la prelucrare zonele în care valorile sunt nedefinite de utilitarul Praat.
- zone unde nu este sunet vocalic;
- zone unde este sunet vocalic, dar utilitarul Praat nu detectează F0 întotdeauna – eroare de detecție;
- în cazul în care pe un sunet vocalic se obțin cu utilitarul Praat valori aberante – mult mai mică sau mai mare decât o valoare în limite normale (de exemplu. 70-500Hz, la bărbați).

Validarea se realizează prin parcurgerea tuturor datelor de intrare din fișierele de adnotare TextGrid create de Praat, iar pe baza unor criterii, sunt eliminate valori eronate furnizate de același utilitar. Fiecărui fișier TextGrid (fișier de adnotare asociat unei înregistrări) îi corespunde un fișier cu valorile formantului F0 precum și un fișier cu valorile formanților F1-F4. Corelarea fișierelor se face în mod automat, absența unuia dintre ele fiind semnalată ca eroare. Aceste criterii sunt stabilite de către utilizator și ele pot fi modificate luând în considerare intervalele de variație ale fundamentalei și a formanților. Segmentele existente în fișierul de adnotare cu markerii respectivi de timp permit validarea sau invalidarea existenței formanților. De exemplu, dacă o pauză este corespunzătoare unui segment între t1 și t2 (markeri de timp), în fișierul de formanți pe același segment se invalidează F0 și formanții F1-F3.

Adăugarea de noi informații se realizează prin reprocesarea datelor. Fișierul binar conține informații despre data creării, versiunea programului WinCollection utilizat, versiunea fișierului binar, despre persoana care a creat colecția, precum și toate informațiile primare. Un alt avantaj al organizării datelor îl constituie posibilitatea reproductibilității rezultatelor obținute.

Cumularea informațiilor într-un singur fișier asigură flexibilitate în accesarea lor, selectarea rapidă a lor pe baza unor criterii de intrare stabilite de utilizator. Nu există date primare, chiar și valorile eronate furnizate de Praat, care să nu se regăsească în colecție.

Structurarea ne permite analize statistice complexe pe un set de date provenit dintr-o selecție a lor și obținerea rapidă a informațiilor asupra variației parametrilor analizați. În această lucrare raportăm numai statistici pentru vocale. Sunetele semivocalice (despre care se știe că au F0 și formanți superiori) nu sunt luate în considerare în rulările pentru această lucrare, rulări în care s-au setat ca zone de interes doar vocalele.

3. Metodologia de procesare a datelor

Programul permite realizarea unei „statistici tomografice”, în sensul că permite determinarea unor „secțiuni de interes” în statisticile realizate. Programul asigură flexibilitate pentru o analiză statistică pe orizontală cât și pe verticală. De exemplu, programul permite analiza unui singur parametru (F0) pentru toți vorbitorii și pentru toate stările emoționale, sau analiza rezultatelor obținute de la un vorbitor în comparație cu toți ceilalți vorbitori. Analiza este organizată la nivelul stărilor emoționale, la nivelul vocalelor și la nivelul valorilor formanților din cadrul vocalelor.

Programul WinCollection are ca date de intrare fișierele .wav, fișierele de adnotare TextGrid, fișierele cu valorile formanților F0 și F1-F3, fișierele Codes.txt, Phrases.txt și Emotions.txt. Fișierul Codes.txt conține date despre toți vorbitorii ale căror înregistrări trebuie să fie în directorul respectiv. Fișierul Phrases.txt este necesar pentru a stoca frazele permise în fișierele de adnotare. Fișierul e utilizat în depistarea erorilor de adnotare. Prin compararea stringurilor din fișierul Phrases.txt cu stringurile corespunzătoare adnotărilor, sunt depistate lipsa de litere, inversiuni de litere sau alte erori de acest tip. F0 din figura 2 notează toate fișierele generate de Praat conținând valorile lui F0 pentru fișierele .wav din directorul respectiv. Similar, fișierele F1-F3 din figura 2 desemnează ansamblul fișierelor generate de Praat conținând acești formanți. Fișierul Emotions.txt include stringul care desemnează emoțiile din înregistrări. Acest fișier este utilizat la verificarea corectitudinii denumirilor directoarelor destinate să conțină fișierele înregistrărilor cu emoții, respectiv și fișierele anexe corespunzătoare. Se observă că mai mult din fișierele de intrare folosite sunt destinate doar etapei de verificare automată a corectitudinii denumirilor și adnotărilor. Fișierele de înregistrare au extensia .wav, fișierele de adnotare au extensia TextGrid. Erorile depistate sunt înscrise într-un fișier special numit Dc.log. Același fișier este folosit pentru înscrierea și a altor tipuri de erori depistate în fazele următoare ale programului. Fișierul va fi util ulterior în corectarea corpusului respectiv.

Programul tratează (blocul tratarea erorilor conform figura 2) și alte tipuri de erori existente în corpus, de exemplu incompletitudinea corpusului (lipsa unor înregistrări cu o emoție, sau existența unui fișier cu valorile F0, dar nu și cu valorile F1-F3, etc.). Erorile depistate sunt incluse în fișierul Dc.log.

Fișierul binar Dc.bin conține toate datele de intrare (mai puțin .wav) cu structuri organizate în forma unui array cu elemente de dimensiuni egale. Principalele elemente sunt structuri - de tipuri diferite. De exemplu, structura vorbitorilor conține informațiile din fișa vorbitorului. Prin această structurare, informația este accesibilă ușor pentru prelucrări statistice. Rezultatele programului sunt fișierele text Dump.txt și Statistics.txt (vezi figura 2).

ASPECTE METODOLOGICE DE ORGANIZARE A DATELOR ȘI DE ANALIZĂ STATISTICĂ A VOCILOR EMOȚIONALE

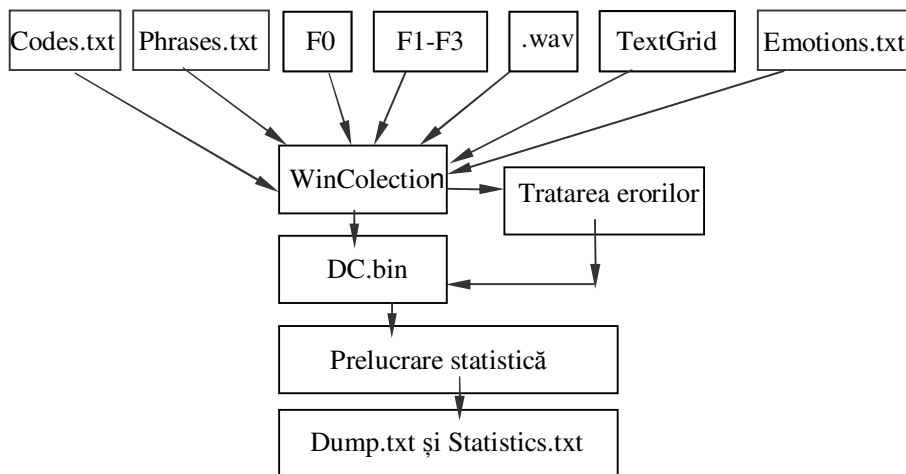


Figura 2: Organigrama programului WinCollection

Programul WinCollection creează colecția de date, fișierul binar DC.bin, care se va obține după parcurgerea și procesarea tuturor fișierelor de intrare. Acest fișier este creat în directorul colecției de date, fiind ulterior folosit de aplicație în analiza statistică a vocilor emoționale.

3.1 Etapele programului

Etapele programului sunt:

1. încărcarea informațiilor despre vorbitori din fișierul Codes.txt;
2. parcurgerea fișierului Phrases.txt și încărcarea informațiilor despre propozițiile pentru care există înregistrări;
3. încărcarea informațiilor despre stările emoționale pentru care există înregistrări, din fișierul Emotions.txt;
4. procesarea datelor de intrare în vederea detectării și semnalării eventualelor erori pe baza informațiilor din fișierele Codes.txt, Phrases.txt și Emotions.txt:
 - 4.1. încărcarea în memorie a datelor de intrare din fișierul de adnotare și a fișierelor cu valorile formanților F0 și F1-F3;
 - 4.2. semnalarea în fișierul Dc.log a pașilor încărcării datelor și semnalarea erorilor;
5. crearea fișierul binar Dc.bin, care va conține întreaga colecție de date;
6. crearea fișierelor text, Dump.txt și Statistics.txt, descrise ulterior.

Programul parcurge toate sub-directoarele care conțin înregistrări despre o anumită stare emoțională, fiecare stare emoțională având în corespondență un sub-director care poartă numele respectivei stări. În fiecare sub-director, programul identifică toate seturile de câte patru fișiere (fișierul .wav, adnotarea corespunzătoare împreună cu fișierele ce conțin valorile formanților F0, F1-F4). Identificarea se face pe baza codului vorbitorului și a propozițiilor pentru care există înregistrări.

3.2 Fișiere de ieșire

Fișierul binar conține și următoarele informații suplimentare: data și ora creării, versiunea de program utilizată, informații despre persoana care a realizat prelucrarea, etc.

Fișierul text Dump.txt conține:

- informații generale despre colecție (data creării, versiunea colecției, versiunea programului WinCollection, etc.);
- un sumar al bazei de date (număr vorbitori, fraze, stări, număr înregistrări pentru fiecare stare și frază);
- datele de intrare;
- un set de rezultate statistice pentru fiecare vocală, inclusiv precizarea numărului de apariții pentru fiecare stare emoțională (rezultatele se referă la minimumul, maximumul, valoarea medie și dispersia, calculate pentru valorile F0 respectiv F1-F3, pentru fiecare vocală, pentru fiecare apariție a vocalei într-o înregistrare).

Fișierul Statistics.txt generat de către programul WinCollection conține următoarele informații:

- pentru fiecare vocală, numărul de apariții pentru fiecare stare emoțională;
- pentru fiecare vocală, valoarea medie pentru populație, pentru formantul F0, respectiv F1-F3;

3.3 Tratarea erorilor

Programul semnalează, în urma prelucrării datelor de intrare, erorile găsite. Principalele erori semnalate de către program în urma prelucrării datelor sunt:

- lipsa adnotării pentru o înregistrare existentă în directorul analizat;
- erori în adnotare (apariția la un anumit nivel a unui fonem care nu apare în propoziția înregistrată);
- absența fișierului care conține informațiile despre formantul F0 pentru o anumită înregistrare;
- absența fișierului care conține informațiile despre formanții F1-F4 pentru o anumită înregistrare.

După corecția erorilor, este necesară re-procesarea tuturor datelor de intrare în vederea obținerii fișierului binar. Nu se fac adăugiri la colecție; adăugarea unei noi înregistrări la o colecție deja existentă se face prin re-procesarea datelor de intrare și obținerea unui nou fișier binar.

În prezent analiza statistică este implementată pe o vocală specifică, pentru un același vorbitor, în cadrul aceleiași propoziții specifice, pentru o stare specificată.

3.4 Estimarea variabilității

Exemplificăm utilizarea programului în vederea obținerii unui coeficient de variabilitate inter și intra-stare pentru valorile formanților precum și o statistică generală a formanților pentru vocale în limba română („a”, „e”, „i”, „u”, „ă”) în funcție de starea emoțională. În cadrul analizei statistice, primul autor a introdus un coeficient de variabilitate (asimetrie) definit conform formulei:

$$\eta = \frac{1}{4} \sum_{k=0}^3 \frac{1}{N} \sum_{\varphi=1}^M \frac{1}{M} \sum_{\psi=1}^M \frac{|F_k[j; \varphi, \psi] - F_k[\psi]|}{F_k[\psi]}, \quad (1)$$

unde: M- este numărul de apariții ale fonemului ψ în fraza φ respectivă, iar k – numărul formantului curent. $F_k[\psi]$ reprezintă valoarea medie a formantului k , pentru fonemul ψ , așa cum este determinată pentru un mare număr de vorbitori, în diverse contexte de pronunție.

Se recunoaște imediat că se realizează media diferențelor absolute ale valorilor formanților respectivi, în pronunțiile date ale fonemului ψ , diferențe normalizate cu valoarea medie a formantului respectiv pentru fonemul respectiv, $F_k[\psi]$:

$$\frac{1}{M} \sum_{\psi=1}^M \frac{|F_k[j; \varphi, \psi] - F_k[\psi]|}{F_k[\psi]}. \text{ Deoarece în sumă apar diferențele unui singur vorbitor față}$$

de toți ceilalți, semnificația este aceea a diferenței vorbitorului față de media vorbitorilor (specificitatea vorbitorului în pronunția fonemului respectiv, în cadrul frazei date, pentru formantul ales). Aceste valori medii devin elemente în următoarea sumă, care, prin multiplicare cu $1/N$ devine media pentru toate frazele (mediile pe frază sunt mediate în sumă după φ). În fine, se face prin ultima sumă media, pentru toți formanții, a diferențelor vorbitor- media vorbitorilor. De remarcat că dacă se calculează numai sumele interioare -care reprezintă rezultate parțiale în program – se determină specificitatea vorbitorului la nivel de formant al unui fonem în fraza specificată, respectiv la nivel de formant al fonemului specificat pe toate frazele. Desigur, pentru fiecare fonem (aici, vocală) și pentru fiecare vorbitor, se obține câte un coeficient de variabilitate.

„Coeficientul de asimetrie” în pronunție a vorbitorului v_j față de media generală la aceeași stare și aceeași vocală, pentru formantul k este dat de formula:

$$\alpha_{B,v,k} = \frac{\overline{F}_{\#k}(\text{starea "B";vocala" a"; vorbitor "v"}) - \overline{F}^0_{\#k}(B,a,\text{toti vorbitorii})}{\overline{F}^0_{\#k}(B,a,\text{toti vorbitorii})}, \quad (2)$$

4. Exemple de rezultate preliminare

Exemplele prezentate nu reprezintă rezultate definitive.

Numărul de ocurențe pentru fiecare vocală în fișierele care au fost analizate sunt date în tabelul 1. Menționăm că vocalele analizate din cele patru propoziții sunt opt „a”, patru „e”, patru „i”, doi „ă” și un „u”, pentru cinci vorbitori, patru stări emoționale cu o medie

de trei rostiri ale fiecărei propoziții. În cazul stării de tristețe am supus analizei doar două propoziții pentru trei vorbitori, iar pentru neutru trei propoziții pentru doi vorbitori.

Tabel 1: Numărul de apariții ale vocalelor analizate în funcție de starea emoțională

Vocale	Starea de bucurie	Starea de furie	Starea de tristețe	Ton neutru
A	84	103	52	102
E	40	50	21	54
I	39	48	24	60
Ă	22	27	14	24
U	11	14	0	9

Exemplificăm, în tabelele 3 și 4, utilitatea metodei expuse pentru caracterizarea unui vorbitor. Dintre cei cinci vorbitori analizați în această lucrare, unul (11861) se distinge de ceilalți prin valorile mult diferite pentru vocalele „u”, „e”, „i”, pentru valorile formaților F0 și F1, iar alt vorbitor (05392) se distinge prin valorile mult mai mici pe vocalele „e” și „a”, pentru valorile F0 și F3. Această concluzie este justificată de valorile prezentate în tabelele 3, 4 și 5. Această distanțare nu se manifestă la formații F1 și F2 în general.

Tabel 2: Valorile coeficientului de variabilitate pentru toate fonemele („a”, „e”, „i”, „u”, „ă”) analizate, pentru toți vorbitorii analizați în funcție de starea emoțională

Codul Vorbitorului	Starea de bucurie	Starea de furie	Starea de tristețe	Tonul neutru
11861	0.156	0.150	0.110	0.178
12121	0.157	0.172	0.102	0.126
83714	0.206	0.166	0.116	0.219
05392	0.093	0.091	0.104	0.059
32167	-	0.102	0.049	0.103

Tabel 3: Valorile coeficientului de asimetrie pentru valorile formaților F0-F3 în funcție de starea emoțională, pentru vorbitor 11861m, vocala „a”

11861/vocala A	Coef. de asimetrie pentru valoarea lui F0	Coef. de asimetrie pentru valoarea lui F1	Coef. de asimetrie pentru valoarea lui F2	Coef. de asimetrie pentru valoarea lui F3
Starea de bucurie	0.33	0.01	0.03	0.003
Starea de tristețe	0.39	0.18	0.28	0.10
Starea de furie	0.27	0.03	0.12	0.01
Ton neutru	0.33	0.10	0.13	0.07

Tabel 4: Valorile coeficientului de asimetrie pentru valorile formaților F0-F3 în funcție de starea emoțională, pentru vorbitor 05392m, vocala „a”

05392/vocala A	Coef. de asimetrie pentru valoarea lui F0	Coef. de asimetrie pentru valoarea lui F1	Coef. de asimetrie pentru valoarea lui F2	Coef. de asimetrie pentru valoarea lui F3
Starea de bucurie	0.11	0.11	0.18	0.07
Starea de tristețe	0.11	0.16	0.12	0.03

ASPECTE METODOLOGICE DE ORGANIZARE A DATELOR ȘI DE ANALIZĂ STATISTICĂ A
VOCILOR EMOȚIONALE

05392/vocala A	Coef. de asimetrie pentru valoarea lui F0	Coef. de asimetrie pentru valoarea lui F1	Coef. de asimetrie pentru valoarea lui F2	Coef. de asimetrie pentru valoarea lui F3
Starea de furie	0.09	0.05	0.12	0.04
Ton neutru	0.08	0.13	0.10	0.02

Tabel 5: Valorile coeficientului de asimetrie pentru valorile formanților F0-F3 în starea de bucurie, pentru patru vorbitori, vocala „u”

Vorbitori/vocala la U/Starea de bucurie	Coef. de asimetrie pentru valoarea lui F0	Coef. de asimetrie pentru valoarea lui F1	Coef. de asimetrie pentru valoarea lui F2	Coef. de asimetrie pentru valoarea lui F3
11861	0.38	0.32	0.31	0.11
12121	-	0.47	0.19	0.16
83714	0.23	0.43	0.29	0.15
05392	0.13	0.24	0.12	0.07

Rezultatele din tabelul 2 arată că valorile coeficienților de variabilitate sunt mai mari în cazul stării de bucurie și sunt mai scăzute în cazul stării de tristețe. În tabelele 3, 4 și 5 valorile coeficientului de asimetrie sunt cele mai mari în cazul frecvenței fundamentale, iar cele mai mici pentru formantul de ordin trei.

S-a urmărit diferențierea unui vorbitor față de ceilalți în funcție de stările emoționale (tabel 2) precum și analiza parametrilor F0-F3 pentru toți vorbitorii, în funcție de vocale (tabel 5). În tabelele 3 și 5 am exemplificat diferențele ce apar între doi vorbitori pentru aceeași vocală „a” pe baza parametrilor F0-F3 în funcție de stările emoționale.

5. Concluzii și direcții viitoare

Scopul principal al acestei lucrări a fost de a prezenta metodologia de organizare a datelor în vederea procesării lor cu o nouă aplicație informatică dezvoltată de colectivul nostru și de a prezenta sumar capabilitățile aplicației. Am evidențiat principalele obiective care au stat la baza metodologiei de organizare și avantajele aplicației realizate. S-au prezentat etapele obținerii colecției de date precum și posibilitatea de manipulare pentru prelucrări statistice. Analiza sumară realizată exemplifică flexibilitatea în utilizare.

În viitor, pe baza informațiilor privind sexul vorbitorului din fișierul Codes.txt se vor adapta limitele acceptabile pentru F0 ce se vor aplica în validarea valorilor din fișierul cu F0, notat F0 în figura 2. În etapa următoare, analiza va fi extinsă la toți vorbitorii din cadrul corpusului SRoL.

Mulțumiri. Cercetarea a fost realizată cu sprijinul Academiei Române, în cadrul temei interne a Institutului de Informatică Teoretică din Iași. Autorii mulțumesc celorlalți co-autori ai sitului Sunetele Limbii Române pentru includerea altor înregistrări noi care au fost analizate în lucrare precum și referențelor care au realizat observații pertinente.

Contribuția autorilor: Primul autor a inițiat tema cercetării și structura lucrării, a elaborat metodologia generală de lucru, conceptul și structura generală a aplicației informatice, a propus formulele menționate

în lucrare și a coordonat întreaga cercetare. Al doilea autor a implementat în C++ programul WinCollections. Al treilea autor a realizat înregistrările, a efectuat manual adnotările și a creat cu utilitarul Praat fișierele necesare utilizării programului WinCollections. Toți autorii au contribuit la redactarea lucrării.

Referințe bibliografice

- Boersma, P., Weenink, D., Institute of Phonetic Science, University of Amsterdam, Praat: doing phonetics by computer, *www.praat.org*.
- Feraru, M., Teodorescu, H.N. (2009). Classification of the Emotional States in Speech using the SRoL Database – Preliminary Results, *Proc. 3rd Int. Conf. Electronics, Computers and Artificial Intelligence*, Pitesti, România, ISBN 1843-2115, 27-32.
- Kienast, M., Sendlmeier, W. F. (2008). Acoustical Analysis of Spectral and Temporal Changes in Emotional Speech, An Acoustic Framework for Detecting Fatigue in Speech Based Human-Computer-Interaction, *Lecture Notes in Computer Science*, ISBN 978-3-540-70539-0, 5105/2008, 54-61.
- Mcgilloway, S., Cowie, R., Douglas-Cowie, E. et al. (2000). Approaching Automatic Recognition of Emotion from Voice: A rough Benchmark. *Proc. ISCA Workshop Speech and Emotion, Newcastle*, 207-212.
- Nakatsu, R., Solomides, A., Tosa, N. (1999). Emotion Recognition and its Application to Computer Agents with Spontaneous Interactive Capabilities. *Proc. IEEE Int. Conf. Multimedia Computing and Systems*, Florence, Italy, 2, 804-808.
- Scherer, K., A. (2000). Cross-Cultural Investigation of Emotion Inferences from Voice and Speech: Implications for Speech Technology. *Proc. Conf. Spoken Language Processing (ICSLP)*, China, http://www.unige.ch/fapse/emotion/publications/pdf/icspl00_crosscul.pdf.
- Teodorescu, H.N., Feraru, M. (2007). A study on Speech with Manifest Emotions, 10th International Conference on Text, Speech and Dialogue, *Lecture Notes in Computer Science, Springer Verlag*, ISBN 978-3-540-74627-0, 4629/2007, 254-262.
- Teodorescu, H.N., Feraru, M., Trandabăț, D., Zbancioc, M., Luca, R., Verbuță, A., M. Ganea, R. Voroneanu, O. Pistol, L. (2005). Proiectul Sunetele Limbii Române, www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.htm.
- Ververidis, D., Kotropoulos C. (2006). Emotional speech recognition: Resources, features, and methods, *Speech Communications*, 48: 9, 1162-1181.

PROGRAM DE EDITARE CONTURI INTONAȚIONALE ÎN LIMBA ROMÂNĂ BAZAT PE IERARHII DE FORME PROSODICE FUNCȚIONALE

VASILE APOPEI, DOINA JITCĂ, OTILIA PĂDURARU

Institutul de Informatică Teoretică

Academia Română - Filiala Iași

vapopei@iit.tuiasi.ro

Rezumat

În lucrare se prezintă un program-editor de arbori intonaționali ce poate fi utilizat în analiza și sinteza intonației în limba română. Editarea unui contur intonațional este legată de un text de intrare (frază) și implică împărțirea acestuia în unități prosodice pe mai multe nivele. Editorul ajută la structurarea unităților în arbori de rostire, precum și la asocierea lor cu funcții și forme de contur F0. Formele aparțin unui set de categorii prosodice funcționale la nivelul actului comunicativ, definite în cadrul modelului de intonație pe care se bazează editorul. În secțiunea 2 se prezintă modelul de intonație, sub aspectul perspectivei ierarhice și funcționale asupra structurii conturului intonațional. În secțiunea 3 se prezintă funcționarea editorului pentru construirea arborilor, pentru gestionarea fișierelor conținând arbori și pentru transformarea arborilor în contururi intonaționale prin sinteză vocală.

1. Introducere

Lucrarea prezintă un program de editare grafică-interactivă a conturilor intonaționale în limba română bazat pe un set de categorii predefinite de forme elementare pentru conturul frecvenței F0. Un astfel de program ar putea fi folosit în învățarea unor elemente noi de modelare a intonației în limba română, cât și în învățarea limbii române ca limbă străină, pentru antrenarea cursanților în formarea unor intonații adecvate diferitelor tipuri de enunțuri. Contururile intonaționale pot fi construite din forme predefinite, funcție de tipul și structura enunțului. Acest tip de program poate fi folosit de specialiștii din domeniul analizei și sintezei vocale dar și în domeniul învățării limbii române, ca limbă străină, fiind un instrument de antrenare pentru producerea pattern-urilor intonaționale în mod conștient. În cadrul unor astfel de activități, studenții vor trebui să învețe semnificația diferitelor pattern-uri intonaționale în cadrul discursului și apoi să le recunoască la nivel perceptual [Rocca, (2007)].

În proiectarea programului s-a folosit experiența acumulată în proiectarea modului de generare a conturului F0, în sistemul text-voce pentru limba română, dezvoltat la Institutul de Informatică Teoretică [Apopei et al., (2007)], [Jitca et al., (2009a)].

Modalitatea de descriere a conturilor intonaționale pe baza unor forme elementare de contur F0, a fost posibilă după identificarea unor categorii funcționale de unități prosodice elementare și a studierii combinațiilor acestora la nivelul unităților melodice neelementare. Definirea categoriilor funcționale s-a bazat pe stabilirea unei relații între conturul F0 al unităților elementare și funcția acestora în realizarea actului comunicativ. În secțiunea 2 se prezintă modelul de intonație, cu aprofundare asupra ierarhiei unităților intonaționale și a realizării conturilor melodice ca secvențe ierarhizate ale

acestor unități. Cunoașterea acestor aspecte este necesară pentru înțelegerea modului de operare cu programul de editare. În secțiunea 3 sunt prezentate funcțiile editorului pentru generarea unei intonații în limba română, precum și managementul structurilor arborescente intonaționale.

2. Modelul de intonație

Modelul intonațional aflat la baza editorului face parte din categoria celor perceptuale ca și modelul IPO [J.'t Hart et.al., (1990)], care idealizează un contur F0 natural prin aproximarea cu o secvență de forme elementare selectate dintr-un set predefinit. Formele de contur F0 corespund unităților prosodice din care se consideră a fi compus conturul F0 analizat. Modelul de intonație pentru limba română, prezentat în [Jitca et al.(2009b)] cuprinde două aspecte importante ce vor fi prezentate în subsecțiunile următoare. Subsecțiunea 2.1 prezintă tipurile de domenii prosodice și realizarea lor prin unități în cadrul arborelui unei rostiri. Subsecțiunea 2.2 prezintă categoriile funcționale în baza cărora se face împărțirea unităților prosodice, astfel încât conturile F0 neelementare să poată fi descrise prin secvențe de unități funcționale.

2.1. Perspectiva ierarhică a structurii conturului intonațional

Pentru a da un înțeles conturilor F0 complexe, s-a procedat la împărțirea acestora în segmente melodice elementare (SME), identificarea pattern-urilor lor melodice și a raporturilor tonale între acestea în cadrul segmentelor melodice neelementare (nSME). SME-urile sunt segmente care includ un singur eveniment fonologic de tip accent și corespund unităților de accentuare (în engl. „*accental unit*”-AU) definite în cadrul modelelor de intonație autosegmental-metrice. La nivel textual, AU-rile corespund unui cuvânt cu accent, sau unui cuvânt însoțit de eventuale cuvinte neaccentuate.

SME-urile prin caracteristica lor tonală (ton țintă sau nivel tonal mediu) intră în contraste tonale cu alte SME-uri din vecinătate sugerând formarea unor grupuri (grupuri de unități de accentuare - în engl. „*accental uni group*” - AUG) cărora le corespund segmentele melodice neelementare (nSME). Noi considerăm că un grup este, în esență, un contrast între doi „poli” tonali și că nu este obligatorie prezența tonurilor de graniță pentru delimitarea lor. Caracteristica tonală a grupului este dată de nivelul țintă dominant sau nivelul tonal mediu, în lipsa unuia țintă dominant. Prin caracteristica tonală dominantă grupurile pot intra în relație de contrast tonal cu elementare vecine ce sunt observabile la un nivel de grupare superior.

Cu ale cuvinte o frază intonațională se poate descompune într-o ierarhie de segmente elementare și neelementare. Spre exemplu, o rostire neutrală a textului *Lui Winston / (i se treziră / (niște vagi /amintiri))* poate fi descrisă prin unitățile demarcate pe text prin „/” și paranteze rotunde iar în ierarhia din figura 1 prin unitățile AU1-AU4 la nivel elementar și AUG1-AUG2 la nivel neelementar.

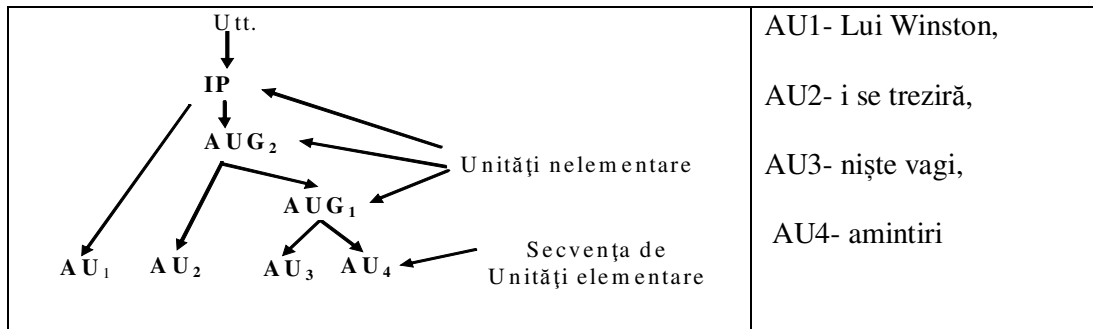


Figura 1: Arborele intonațional al unei rostiri neutrale a textului „Lui Winston / (i se treziră / (niște vagi /amintiri))”

Modalitatea prin care se proiectează diferitele nivele ale unităților din cadrul arborelui intonațional în cadrul rostirii, fie în varianta naturală sau sintetizată, o constituie gama de variație a conturului F0 în cadrul acestora (unităților de nivel mai înalt le corespund game mai largi ale variației frecvenței F0). Nivelele arborelui din figura 1 sunt mai bine puse în evidență de reprezentarea echivalentă din figura 2. Pe frunzele arborelui sunt unitățile de tip AU iar în noduri cele de tip grup de AU (AUG). Ierarhia modelului poate conține un nivel corespunzător frazei intonaționale (Nivel I), două nivele de AUG (Nivel II și Nivel III) iar pe ultimul nivel sunt situate unitățile prosodice elementare (AU).

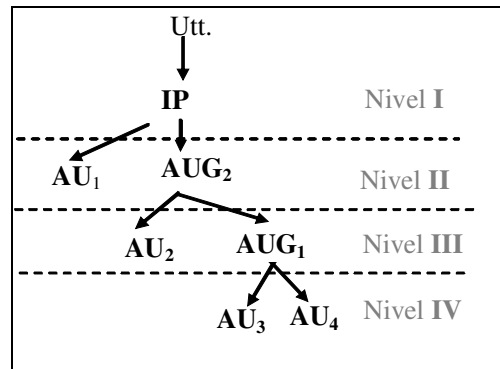


Figura 2: Arborele intonațional al rostirii configurat ca secvențe funcționale pe mai multe nivele

Conturul F0 este redat în sinteză prin concatenarea formelor de contur F0 ale unităților elementare (AU) după poziționarea lor în spațiul (timp, frecvență F0), determinată atât de structura arborelui intonațional cât și de funcțiile acestora în cadrul grupurilor din care fac parte.

2.2. Perspectiva funcțională asupra conturului F0

Pattern-urile conturilor intonaționale elementare au fost inventariate în urma analizei unui corpus de voce în limba română și au fost împărțite pe categorii funcționale la nivelul actului comunicativ. Categoriile au fost asociate unor etichete pentru a fi folosite în descrierea conturilor F0 [Jitca et. al.(2008)]. Perspectiva modelului nostru asupra intonației echivalează din punct de vedere funcțional, unitățile elementare cu cele

neelementare, și ca urmare, funcțiile sunt aceleași cu cele de la nivelul unităților elementare, prezentate în tabelul 1.

Tabelul 1: Semnificația etichetelor funcționale aplicate unităților intonaționale

Eticheta	Funcție comunicativă	Etichetă	Funcție comunicativă
PH	<u>PUSH</u> – „împinge” înainte actul de comunicare prin tonuri înalte care atrag atenția ascultătorului	PH+f	Realizează un eveniment PUSH și o focalizare neutrală minoră
PO	<u>POP</u> – relaxează actul comunicativ fie pentru închiderea unității (caz afirmativ) fie pentru începerea unităților din finalul interogațiilor de tip Da/NU	PH+F	Realizează un eveniment PUSH și o focalizare neutrală majoră
PD	<u>PUSH-DOWN</u> – efectuează un eveniment de tip PUSH și apoi aduce tonul la nivel jos (sil. Neacc.)	PD+F	efectuează un eveniment de tip PUSH și apoi o focalizare neneutrală (sil. acc.)
PU	<u>POP-UP</u> . - efectuează un eveniment de tip POP și apoi ridică tonul sugerând continuarea	PO+f	Realizează un eveniment POP și o focalizare neutrală minoră
f	<i>Focus minor</i> - realizează focalizare neutrală minoră (cuvintele funcționale la nivel gramatical)	PO+F	Realizează un eveniment POP și o focalizare neutrală majoră
F	<i>Focus major</i> -realizează focalizare neutrală majoră (cuvintele nefuncționale la nivel gramatical)	L	<u>LINK</u> – realizează o funcție de legătură între două unități funcționale contrastante

Aceste categorii de funcții sunt implementate la nivelul modulului de sinteză vocală prin seturi de prototipuri de contururi F0 elementare, care generează la nivelul percepției segmente melodice elementare.

3. Prezentarea programului de editare

Programul facilitează construirea structurilor intonaționale ierarhice ale rostirilor (arborii rostirilor – în engleză, *utterance tree*) prin manevrarea unor obiecte vizuale corespunzătoare segmentelor elementare și neelementare, până la realizarea unui contur intonațional dorit corespunzător rostirii unui text în limba română. Arborii de rostiri constituie intrările/ieșirile unui astfel de editor asupra cărora se pot aplica următoarele operații:

- **Creare** - construire arbore nou
- **Editare** – modificarea arborelui curent dintr-un fișier XML
- **Salvare** - pentru memorarea arborelui editat într-un fișier text de tip XML, prin înlocuirea celui inițial sau prin adăugare.
- **Vizualizare** – pentru afișarea grafică a arborilor intonaționali conținuți într-un fișier XML.
- **Redare** – pentru a comanda sinteza vocală a rostirii pe baza a arborelui intonațional curent și a textului conținut de acesta.

Pentru operația de **salvare** a arborilor generați prin comenzile **creare/editare** s-au folosit fișiere XML în care textul corespunzător arborelui intonațional este adnotat cu un set de tag-uri și atribute prezentate în tabelul 2.

PROGRAM DE EDITARE DE CONTURI INTONAȚIONALE ÎN LIMBA ROMÂNĂ BAZAT PE IERARHII DE FORME PROSODICE FUNCȚIONALE

Tabelul 2: Tag-uri și atribute pentru adnotarea prosodică a unui text

Tag	Atribut	Valoare	Tipul unității
<IP>			Frază intonațională
<AUG>	<i>Function</i>	PH, PO, PU, PD, PH+f, PO+f, PD+F, PH+F, PO+F,f , F, L	Grup de unități de accentuare
	<i>TonalContrast</i>	s(trong),w(eak)	
<AU>	<i>Function</i>	PH, PO, PU, PD, PH+f, PO+f, PD+F, PH+F, PO+F,f , F, L	Unitate de accentuare
	<i>PichAccent</i>	H*, L*, L+H*, H+!H*, L*+H	
	<i>TonalContrast</i>	s(trong),w(eak)	
<Syl>	<i>Length</i>	Small, Medium, Large	Silabă
	<i>Energy</i>	Small, Medium, Large	

Folosind tag-urile, cu atributele și valorile din tabelul 2, arborii intonaționali creați de editor, ca ierarhii de obiecte, corespunzătoare unităților prozodice, ce includ textul aferent, se transformă în urma comenzii de **Salvare** în structuri XML de text adnotat prozodic.

Operația de **Creare/Editare** a arborilor implică definirea în memorie a unor obiecte de tip „segment contur intonațional” corespunzătoare unităților prozodice (unități de accentuare și grupări ale acestora). Aceste obiecte au câte un echivalent vizual în fereastra principală, de editare arbore, ce face posibilă manevrarea lor cu mouse-ul (figura3).

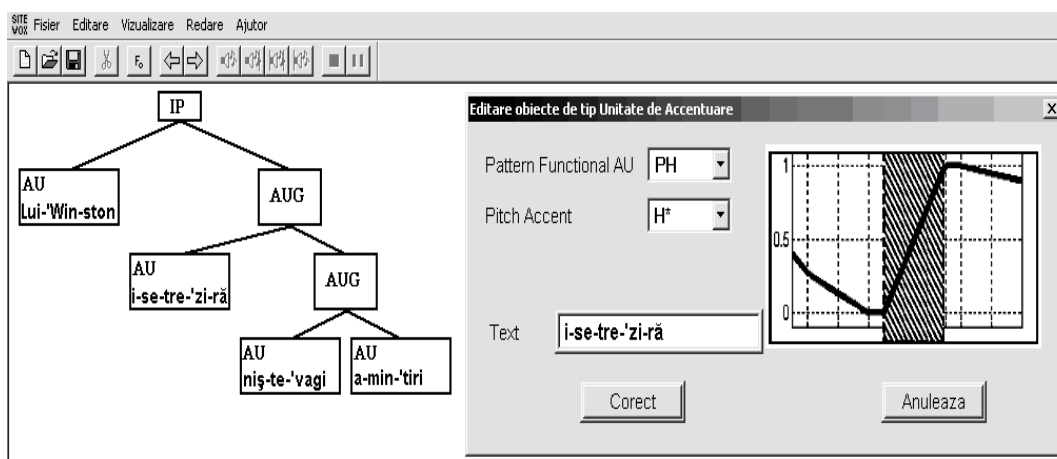


Figura 3: Fereastră de editare a unui arbore intonațional

Comenzile de bază aplicate obiectelor vizuale pentru construirea unui arbore sunt:

- **Selectie** – selecție obiect/grup obiecte
- **Inserare** – inserare unitate AU
- **Ștergere** – ștergere Unitate AU/AUG
- **Grupare** – grupare unități AU/AUG selectate în vederea creării unui nou AUG

Comenzile de inserare, ștergere se aplică în raport cu obiectul selectat. Comanda de grupare se aplică obiectelor adiacente selectate, aparținând aceluiași nivel.

Fiecare obiect are un set de atribute dintre care funcția unității la nivelul actului comunicativ este cel mai important pentru că acesta este folosit în selectarea pattern-ului de contur F0 și/sau poziționarea acestuia în gama frecvenței F0. Aplicarea valorilor pentru atribute și editarea textului aferent unităților de accentuare se face cu comanda **Editare** în cadrul unei ferestre de dialog - figura 3.

4. Concluzii

Programul de editare de arbori intonaționali este un instrument util pentru generarea de intonații corecte în limba română corespunzătoare diferitelor structuri gramaticale și sintactice. Pentru construirea unui arbore care să genereze o intonație validă operatorul trebuie să cunoască modelul intonațional bazat pe ierarhii de forme prosodice funcționale. Editorul va putea fi folosit în proiectarea de reguli pentru modulul de predicție prozodică. În cazul folosirii metodelor de învățare automată se pot construi baze paralele de arbori intonaționali și sintactici pe baza cărora să se poată deduce reguli de punere în corespondență.

Odată cu implementarea unui model de intonație pentru o anumită limbă a apărut ca o aplicație imediată instruirea studenților străini în învățarea intonației limbii respective. Dar instrumentul poate fi folosit chiar și pentru a instrui vorbitorii nativi ai limbii respective în a gândi și înțelege un contur intonațional prin prisma elementelor modelului.

Referințe bibliografice

- Apopei, V., Jitcă, D. (2007). Module for generating the F0 Contour using as input a Text structured by prosodic information, *Advances in Spoken Language Technology, Romanian Academy*, 119-126.
- Rocca, A. P. D. (2007). New Trends on the Teaching of Intonation of Foreign Languages, *Proceedings of the Fifth International Symposium on the Acquisition of Second Language Speech*, 420-428.
- J.'t Hart, Collier R., Cohen A. (2006). A perceptual study of intonation, Cambridge University Press.
- Jitcă, D., Apopei V. (2009a). A Prosodic Control module for a Romanian TtS System, based on melodic contour dictionaries, *Advances in Spoken Language Technology, Romanian Academy*, 77-86.
- Jitcă, D., Apopei, V. and Jitcă, M., (2009b). The F0 Contour Modelling as Functional Accentual Unit Sequences, *International Journal of Speech Technology*, DOI 10.1007/s10772-009-9055-3, Nov, 2009.
- Jitcă, D., Apopei V. (2007). Corpus de voce pentru limba română adnotat cu etichete funcționale la nivelul unităților de accentuare, *Lucrările atelierului „Resurse lingvistice și instrumente pentru prelucrarea limbii române*, p.31-39, Iași, 15.11.2007.

CORPUS PENTRU GNATOFONIE: PROTOCOL, METODOLOGIE, ADNOTARE

HORIA-NICOLAI TEODORESCU^{1,2}, ALINA UNTU¹

¹*Universitatea Tehnică „Gheorghe Asachi”, Facultatea de Electronică,
Telecomunicații și Tehnologia Informației, Iași - România*

²*Institutul de Informatică Teoretică, Academia Română, Iași - România;*

{hteodor, auntu}@etti.tuiasi.ro

Rezumat

Prezentăm elemente privind extinderea unui micro-corpuz de sunete gnatofonice, exemple de analiză formantică și temporală a consoanelor fricative în cazuri de protezare și edentație, precum și rezultatele preliminare.

1. Introducere

Gnatofonia a fost introdusă de către primul autor, parțial la sugestia Prof. Leonid Teodorescu și reprezintă o metodă de analiză a deficiențelor de pronunție a sunetelor datorate diverselor patologii ale aparatului stomatognat (edentație, malocluzie, disfuncții ale articulației temporo-mandibulare, etc.). Primul micro-corpuz de sunete gnatofonice (înregistrate pe casetă) a fost realizat de către primul autor în colaborare cu C. Morărașu, V. Burlui și D. Leca și a stat la baza primelor validări semiempirice în gnatofonie, raportate intern în teza de absolvire a C. Morărașu, iar apoi, cu unele dezvoltări datorate în principal lui C. Morărașu, într-un set de lucrări comunicate sau publicate în anii 1990 (Burlui, Teodorescu, Morărașu, 1993 a, b, 1994, 1996). Un alt corpuz care a fost introdus pe sit-ul „Sunetele limbii române” la secțiunea „Arhiva pentru aplicații de gnatosonie și gnatofonie”, descris în (Teodorescu et al., 2005-2007), (Teodorescu et al., 2007), (Teodorescu, Feraru, 2007 a) a fost realizat de primul autor în colaborare cu S.M. Feraru și a fost inclus în (Teodorescu, Feraru, 2007 b). Dacă termenul de gnatofonie și conceptul de disciplină au fost introduse în anii 1980, observarea empirică a modificărilor fonatorii induse de edentație sunt cunoscute încă din anul 1959 când s-au realizat primele studii sistematice de către Rathbone și Snidecor (Rathbone, Snidecor, 1959). Autorii au studiat instrumental diferențele ce apar în pronunția de propoziții test ce conțin consoane post-dentale (*n, t, d, r, l, s, sh, z, zh, y*), linguo-dentale (*th* sonor și nesonor), labio-dentale (*f, v*) și combinații post-dentale (*ch, j*) înainte și după tratamentul ortodontic a 10 pacienți cu diferite stadii de malocluzie. Au constatat că înainte de tratament din 16 sunete dentale testate un procent de 6.4 % sunt defectuoase, cele mai mari erori concentrându-se asupra consoanelor dentale de tip fricativ (*s, z, sh, zh și th*). După patru ani, opt pacienți din 10 au fost reexaminați după aceeași procedură, rezultând o scădere a procentului de sunete cu pronunție defectuoasă de până la 1.5 %. Erorile de pronunție s-au păstrat la consoanele fricative având însă un grad mai scăzut, iar celelalte sunete dentale au fost corectate astfel încât erorile au fost indetectabile. Acest studiu a dovedit că defectele de vorbire se corectează după tratamentul ortodontic fără a fi nevoie de aplicarea unei terapii de educare a vorbirii. Analiza de semnal vocal în cazul unor patologii ale aparatului stomatognat este importantă în diagnosticarea și evaluarea eficienței tratamentului ortodontic.

Spre sfârșitul anilor 1990, bazându-se pe un număr de 42 de articole consultate, Johnson și Sandy (Johnson, Sandy, 1999) au prezentat stadiul cercetărilor realizate până în anul 1999 de către specialiști în domeniu, vizavi de relația dintre poziția dinților și vorbire. S-a constatat că dinții au un rol important în vorbire, însă relația dintre poziția dinților și vorbire a rămas încă controversată. Se cunoaște capacitatea pacienților de a-și compensa defectele de vorbire ce apar în cazul malocluziilor, însă mecanismul ce stă la bază rămâne necunoscut. Deși s-a demonstrat că există o relație între defectele de dentiție și defectele de pronunție, nu s-a putut concluziona că ar exista o corelație și cu severitatea malocluziilor.

2. Metodologia de culegere a datelor

Metodologia, protocolul de înregistrare și protocolul de documentare au fost preluate de pe sit-ul „Sunetele limbii române”. Subiecții au fost informați anterior înregistrărilor de obiectivele proiectului, fiind asigurați de confidențialitatea datelor personale. Subiecții au semnat un consimțământ informat în conformitate cu protocolul de protecție a subiecților umani și cu principiile etice ale cercetărilor care implică ființa umană existente la nivel național și internațional. Cercetarea pe subiecți umani privind analiza sunetelor gnatofonice și gnatosonice a fost aprobată de către consiliul Facultății de Electronică, Telecomunicații și Tehnologia Informației din cadrul Universității Tehnice „Gheorghe Asachi” din Iași (Teodorescu et al., 2005-2007)¹.

2.1 Metodologia de înregistrare

Sunetele gnatofonice au fost înregistrate cu ajutorul unui microfon prevăzut cu căști, A4 Tech Stereo HS-60, având caracteristicile: frecvență de răspuns 20 Hz-20 kHz, impedanță: $U=3V$, $R=1,5k\Omega$, sensibilitate: $-58dB\pm 2$. Placa de bază a calculatorului pe care au fost efectuate înregistrările este Sony Vaio MBX-189 având încorporată o placă de sunet Intel® High Definition Audio compatible 3D audio (Direct Sound 3D support) cu următoarele caracteristici: procesor de semnal 44-kHz / 16-bit stereo CD quality, mod de ieșire a sunetului 8 canale, 192 kHz / 32 bit, standard Intel HD audio.

Pentru înregistrări s-a utilizat programul GoldWave™, versiunea 5.54, la o frecvență de eșantionare de 22050 Hz cu atributele PCM signed (16-24 bits mono). Culegerea de semnal vocal se realizează în condiții de zgomot redus (amplitudinea zgomotului trebuie să fie mai mică cu cel puțin 20 dB decât amplitudinea frecvenței fundamentale). Conform metodologiei de înregistrare de pe situl SRoL, se recomandă ca poziția microfonului să fie mai jos de gură, aproximativ în dreptul bărbiei (la câțiva centimetri de aceasta), iar distanța de la bărbie să fie aproximativ egală cu distanța până la buze (Teodorescu et al., 2005-2007).

2.2 Metodologia de documentare

Fișele subiecților „Profil vorbitor”, „Chestionar Patologie Vocală și Factori Obiectivi” și „Fișă dinți subiect” au fost completate de către al doilea autor. Pentru „Fișă dinți subiect” s-au luat în considerare informații ce pot fi cunoscute de către subiectul

¹ Avizul este accesibil la:

www.etc.tuiasi.ro/sibm/romanian_spoken_language/images/jpg/doc_subiecti_umani_analiza_gnatofonie.jpg

înregistrat sau vizualizate de către persoana care completează fișa, de tip prezență / absență dinte, plombă, implant, coroană, punte. Am considerat elemente de interes și plombele întrucât acestea dacă nu sunt realizate corect pot modifica structura dintelui având repercusiuni asupra spațiului interdental, în consecință și asupra vorbirii (spațiile dintre incisivi pot produce în vorbire efectul de „fonfăit”). Pentru a valida aceste idei este necesară o analiză statistică pertinentă (cel puțin 20 de subiecți). Un exemplu de fișă ce conține informații referitoare la patologiiile arcadei superioară și inferioară este dat în figura 1. Subiectul prezintă opt plombe simple la nivelul maxilarului și opt plombe simple la nivelul mandibulei, care afectează primii molari de pe fiecare hemiarcadă.

Fișă dinți vorbitor															
Cod: 1234															
Hemiarcada dreaptă superioară								Hemiarcada stângă superioară							
1.8	1.7	1.6	1.5	1.4	1.3	1.2	1.1	2.1	2.2	2.3	2.4	2.5	2.6	2.7	2.8
	Pb	Pb	Pb	Pb							Pb	Pb	Pb	Pb	
4.8	4.7	4.6	4.5	4.4	4.3	4.2	4.1	3.1	3.2	3.3	3.4	3.5	3.6	3.7	3.8
	Pb	Pb	Pb	Pb							Pb	Pb	Pb	Pb	
Hemiarcada dreaptă inferioară								Hemiarcada stângă inferioară							
Legendă: Dinte lipsă - X Plombă - Pb Coroană - C Implant - I Punte - Pu															

Figura 1: Exemplu de fișă dinți subiect.

2.3 Listă cuvinte

Pentru realizarea corpusului de sunete gnatofonice s-a utilizat setul de cuvinte stabilit anterior de primul autor (©), cuvinte ce conțin consoanele *s, f, ș, v, z, j*. S-au ales aceste cuvinte, întrucât studiile realizate până în prezent de către cercetători în domeniu au arătat că fricativele prezintă mai frecvent modificări de pronunție în cazul existenței unor patologii ale aparatului stomatognat (Rathbone, Snidecor, 1959). Analiza comparativă a pronunției consoanelor *s, f, ș*, și a celor semi-vocalice *v, z, j*, furnizează informații despre afectarea dentiției.

Cuvintele utilizate pentru înregistrările gnatofonice sunt: vată / fată, var / far, vuiet (pronunțat vvvvuiet) / vuiet (pronunțat normal, scurt, vuiet) / fffffui / fui, vvvvvalet / valet / fffffailet / failet, vecin / fecior, vvvvvânt / vânt / fffffân / fân, vvvvvine / vvvine / fffffine, vine / fine, vehement / ferment, vierme / fierbe, vâjâit / gâjâit / sâsâit / fâsâit, bâzâie / zâzâie, bâzzzzzâie / zâzzzzzâie, fâșâit / fâlfâit, vâjâie / fâlfâie / sâsâie / fâșâie, vâjjjjjâie / fâlfâie / sâsssssâie / fâșșșșșâie, vâjjjjjâit / sâsssssâit / fâsssssâit / fâșșșșșâit / fâlfâie, vorbit / fortuit / sortit, suit / vuit.

Scopul cercetării noastre este de a identifica și analiza modificările ce apar în dinamica formașilor / pseudo-formașilor consoanelor fricative pronunțate de către două clase de vorbitori: o clasă martor formată din vorbitori cu dentiție normală (neafectată de patologii ale aparatului stomatognat) și o clasă de vorbitori cu defecte de dentiție. Setul de cuvinte întocmit în acest scop include cuvinte ce conțin consoanele *v, f, s, ș, j, z* în context consoană-vocală (CV) sau vocală-consoană-vocală (VCV). Termenul de „pseudo-formași” a fost folosit de primul autor în (Teodorescu, Feraru, 2008), pentru a desemna traseele asemănătoare cu cele ale formașilor pe care le detectează utilitarul

PraatTM în cazul consoanelor fricative. Sunetele neperiodice cum ar fi siflantele *s, f, ș*, nu prezintă frecvență fundamentală, drept urmare nu au nici formanți propriu-ziși (care au, prin definiție, frecvențe egale cu multipli ai frecvenței fundamentale).

Pronunțiile celor 56 de cuvinte pentru fiecare din cei 10 vorbitori sunt înregistrate într-un singur fișier .wav care este salvat cu un nume mnemonic, în formatul „cod_subiect_sex” (ex. 1234_m). Fișierele au fost filtrate și re-salvate în format .wav și în format .ogg (16 și 24 biți). Înregistrările au fost ascultate pentru o analiză perceptuală a zgomotului, iar pronunțiile cu trunchiere (saturație) în amplitudine au fost eliminate.

3. Discuția corpusului

Corpusul de sunete gnatofonice este disponibil pe sit-ul „Sunetele limbii române”, în cadrul arhivei pentru aplicații de gnatofonie și gnatosonie (Figura 2). În paralel am realizat și un corpus de sunete gnatosonice pe care îl vom prezenta într-o lucrare ulterioară.

Arhiva pentru aplicații de gnatofonie și gnatosonie

Metodologia de culegere a semnalelor gnatofonice este identică cu cea de culegere de semnal vocal.

Cuvintele utilizate pentru înregistrările gnatofonice sunt alese astfel încât să se poată analiza comparativ modificările de siflante, fricative și de consoane semi-vocalice.

Cuvintele utilizate pentru înregistrările gnatofonice sunt: *vată / fată; var / far; vuiet (pronunțat wwwuiet) / vuiet (pronunțat normal, scurt, vuiet) / fui / vaiet (pronunțat wwwaiet) / vaiet (pronunțat vaiet) / faieton / vecin / fecior / vânt (pronunțat wwwânt) / vânt (pronunțat vânt) / fân / wwwine, wvine, wvine / vine / fine / vehement / ferment / vierme / fierbe / vâjâit / wwwâjâit / wwwâjâie / fffâșșșșâie / fffâșșșșâit / fâșâit / sâsâit / sssssâssssâie / gâjâit / zâzâie / bââzzzzâââie / bâz&226;ie.*

Un numar de 10 înregistrări gnatofonice provenite de la 5 subiecți de gen feminin și 5 subiecți de gen masculin sunt adnotate, iar fișierele TextGrid și wav corespunzătoare se pot vizualiza (accesa) la adresa (http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/ro/sunete_gnatofonice.htm) sau aici.

Fișă vorbitor	Gnatofonie		Gnatofonie - filtrate		Gnatosonie		Gnatosonie - filtrate	
	wav	ogg	wav	ogg	wav	ogg	wav	ogg
Fișă vorbitor #394715	394715_f	394715_f	394715_f	394715_f				
	394715_f_v	394715_f_v	394715_f_v	394715_f_v1				
	394715_f_v2	394715_f_v2	394715_f_v2	394715_f_v2				
Fișă vorbitor #231515	231515_m	231515_m	231515_m	231515_m				
Fișă vorbitor #280269	280269_m	280269_m	280269_m	280269_m				
Fișă vorbitor #20048					20048_f_v0	20048_f_v0	20048_f_v0	20048_f_v01
Fișă vorbitor #55555	55555_f	55555_f	55555_f	55555_f	55555_f	55555_f	55555_f	55555_f
	55555_f_v1	55555_f_v1	55555_f_v1	55555_f_v1				
Fișă vorbitor #1234	1234_m	1234_m	1234_m	1234_m	1234_m	1234_m	1234_m	1234_m
Fișă vorbitor #1357	1357_m	1357_m	1357_m	1357_m	1357_m	1357_m	1357_m	1357_m
Fișă vorbitor #2001	2001_f	2001_f	2001_f	2001_f	2001_f	2001_f	2001_f	2001_f

Figura 2: Arhiva pentru aplicații de gantofonie și gnatosonie.

3.1 Statistica înregistrărilor și patologiilor

Am realizat înregistrări de sunete gnatofonice provenind de la cinci subiecți de gen masculin și cinci subiecți de gen feminin, cu vârste cuprinse între 21 și 46 ani,

CORPUS PENTRU GNATOFONIE: PROTOCOL, METODOLOGIE, ADNOTARE

majoritatea cu studii superioare, fără afecțiuni respiratorii, laringeale, neurologice, sau psihologice.

Corpusul realizat are la bază înregistrări provenind de la nouă subiecți fără defecte majore de dentiție și fără deficiențe de vorbire detectabile prin percepție auditivă și un subiect cu edentație majoră (13 molari lipsă). Starea dentiției subiecților înregistrați este ilustrată în Tabelul 1.

Tabel 1: Starea dentiției subiecților înregistrați

Cod subiect	Sex	Vârș-tă (ani)	Nr. Pb (plombă) / localizare	Nr. X (lipsă dinte) / localizare	Nr. C (coroană) / localizare	Nr. Pu (punte) / localizare	Tip voce / patologii	
2404	M	30	1 primul incisiv HSS	-	-	-	Profesională, fără patologii	
2001	F	26	7 primul molar HDI, al 2-lea și al 3-lea molar HDS și HSS, al 2-lea molar HSI, al 3-lea molar HDI	-	-	1 al 2-lea molar HDI	Neprofesională, fără patologii	
4312	F	29	3 primii incisivi HSS și HDS, al 3-lea molar HDS	1	al 3-lea molar HDI	2 al 2-lea molar, HSS, al 3-lea molar HSI	Profesională, fără patologii	
1357	M	30	2 canin HSS, al 2-lea molar HSS	2	ultimul molar HDI și HSI	2 al doilea molar HDI și HSI	1 molarii 1, 3, HDS	Neprofesională, fără patologii
2202	F	26	3 al 2-lea molar HSS, al 3-lea molar HDS, al 4-lea molar HSI	1	al 3-lea molar HSI	-	-	Neprofesională, fără patologii
1234	M	30	16 primii 4 molari de pe fiecare hemiarcadă	-	-	-	-	Neprofesională, fără patologii
3371	F	21	4 primul incisiv HSS, al 2-lea molar HSS, al 4-lea molar HDI, al 3-lea molar HSI	2	al 3-lea molar HDI și HSI	-	-	Neprofesională, fără patologii
3298	M	30	3 2 pe primii incisivi de pe HSS și HDS, al 3-lea molar HDS	-	-	2 al 2-lea molar HDI și HSI	-	Profesională, fără patologii
01321	F	30	2 canin și primul molar HDI	-	-	1 primul molar HSS	-	Neprofesională, fără patologii
5343	M	46	1 ultimul molar HSI	13	incisivii 1, 2 HSS, primii 3 molari HDS, ultimii 3 molari HSS, molarii 2,3,4 HSI, molarii 3, 4 HDI	-	-	Neprofesională, fără patologii

În studiul nostru ne interesează în special afectarea incisivilor, caninilor și primilor molari deoarece aceștia intervin în articularea sunetelor. Cei patru subiecți la care am efectuat analiza consoanelor fricative (1234m, 5343m, 01321f, 3371f) prezintă

următoarele patologii relevante: primul subiect prezintă câte o plombă la fiecare din cei patru molari primari; subiectul al doilea prezintă edentație ce interesează primul și al doilea incisiv de pe hemiarcada stângă superioară (HSS) și primul molar de pe hemiarcada dreaptă superioară (HDS); al treilea subiect prezintă câte o plombă pe caninul și primul molar de pe hemiarcada dreaptă inferioară (HDI) și o coroană pe primul molar de pe HSS; cel de-al patrulea subiect prezintă o plombă pe primul incisiv de pe HSS.

4. Metodologia de analiză

4.1 Metoda de prefiltrare

Toate fișierele au fost prefiltrate cu un filtru trece bandă cu frecvențele de tăiere de cca. 70 Hz și 7 kHz și cu atenuare de 100 dB la frecvențele de 50 Hz și 12 kHz (în afara benzii de trecere), setat în utilitarul GoldWave™ (Figura 3). Filtrul se poate aplica o dată sau, repetat, de două ori (echivalent, două filtre înseriate) astfel ca nivelul final de zgomot la 50 Hz să fie cu cel puțin 20 dB mai mic decât nivelul la frecvența fundamentală (F0). În cazul înregistrărilor curente, filtrul s-a aplicat o singură dată.

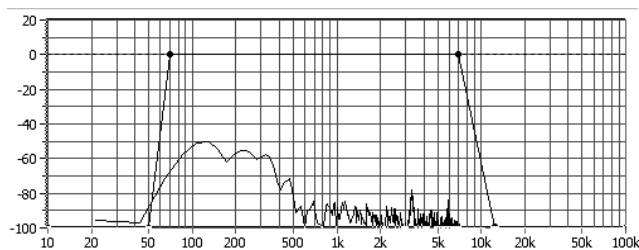


Figura 3: Filtru trece bandă de la 70 Hz la 7 kHz setat în utilitarul GoldWave™.

4.2 Metodologia de adnotare

Pentru o analiză ulterioară a diferențelor / similitudinilor ce apar la nivelul consoanelor fricative, am segmentat și adnotat cele 10 fișiere la nivel de fonem / silabă / cuvânt cu ajutorul utilitarului Praat™ (Boersma, 2002). Am realizat împărțirea fișierelor în subfișiere ce conțin consoanele *f*, *s*, *ʃ*, *v*, *z*, *j*, în context CV și VCV. Au rezultat câte nouă sau 11 subfișiere la fiecare subiect, numite astfel: *vf_CV(a)_cod_subiect_sex* – consoanele *v* și *f* în context CV, *V=vocala a*, *jsfshz_VCV(a-)_cod_subiect_sex* – consoanele *j*, *s*, *f*, *ʃ*, *z* în context VCV, *V=vocala â*, *vfs_CV(o)_cod_subiect_sex* – consoanele *v*, *f*, *s* în context CV, *V=vocala o*, etc. Pentru *ʃ*, *â* și *ă* am utilizat notațiile *sh*, *a-* și *a+*.

În cadrul procesului de adnotare am luat în considerare și pauzele intravorbire (ce apar în rostirea de silabă, cuvânt) notate cu simbolul \$, pauzele intervorbire (dintre cuvinte) fără notație (blanc) și pauzele care nu se aud, notate cu simbolul %. Adnotarea s-a realizat ținând cont de percepția auditivă, prezența / absența frecvenței fundamentale (ce evidențiază caracterul vocalic, sonor, cu activarea corzilor vocale), forma de undă, energia și modificările ce apar pe spectrogramă de la un fonem la altul.

4.3 Metodologia de analiză formantică și temporală

Cu ajutorul utilitarului Praat™ am extras F0, formații / pseudo-formații F1, F2, F3, F4 și timpii corespunzători fiecărui fonem și fiecărui cuvânt în parte. Pentru extragerea formațiilor / pseudo-formațiilor F1, F2, F3, F4, Praat-ul utilizează o fereastră glisantă de analiză de lungime 0,025 s, cu un pas de deplasare de 0,00625 s, preaccentuare de 50 Hz și frecvențe maxime pentru formații de 5500 Hz pentru voci feminine și 5000 Hz pentru voci masculine. Conform manualului de utilizare al utilitarului Praat™, preaccentuarea de 50 Hz determină nemodificarea spectrului la frecvențe sub 50 Hz și amplificarea cu 6 dB la fiecare octavă: 6 dB la frecvența de 100 Hz, 12 dB la frecvența de 200 Hz, etc. (Boersma, 2002).

Pentru extragerea frecvenței fundamentale (F0), manualul Praat™ recomandă utilizarea benzilor de 75 Hz-300 Hz pentru voci masculine, respectiv 100 Hz-500 Hz pentru voci feminine, precum și a unui pas de deplasare al ferestrei de analiză de 0,01 s. „Se folosesc intervale standard pentru F0 la voci masculine și feminine întrucât alegerea unei valori prea mici determină pierderea modificărilor rapide ale lui F0, iar alegerea unei valori prea mari duce la pierderea valorilor foarte mici ale lui F0” (Boersma, 2002).

Analiza a constat din următoarele etape: i) Segmentarea, în vederea analizei modificărilor duratelor fonemelor. Segmentarea este efectuată la nivelurile fonem, silabă și cuvânt, pentru a se determina momentele de timp corespunzătoare granițelor dintre foneme și cuvinte. Într-o fază ulterioară, vom realiza o segmentare mai fină, care să permită determinarea duratelor *glissando*-urilor fonemelor (de ex., *va, fa, rfff* – de ex. în „*până-n vârffful muntelui*” etc.); ii) Generarea fișierelor .txt cu instrumentul Praat, cu valorile formațiilor în evoluția lor; iii) Importarea fișierelor într-un utilitar (de ex., Excel™) și separarea manuală a fonemelor în funcție de limitele temporale corespunzătoare fiecărui fonem; iv) Analiza dinamicii și analiza statistică a formațiilor pe segmentele de interes, comparativ între subiecții cu afecțiuni și tratamente ale aparatului stomatognat și subiecții sănătoși. Am analizat dinamica formațiilor / pseudo-formațiilor consoanelor *v, f, s* din cuvintele *vată, fată, sâsâit, fâsâit, sâsssâit, fâsssâit*. Acestea s-au aplicat la patru vorbitori (doi de gen feminin și doi de gen masculin).

5. Exemple de rezultate preliminare

Studiul realizat pe baza corpusului de sunete gnatofonice s-a axat pe analiza defectelor de pronunție cauzate de patologii ale aparatului stomatognat. S-a urmărit efectuarea unei comparații între dinamica formațiilor consoanelor fricative la subiecți cu dentiție normală și subiecți cu defecte de dentiție.

În continuare prezentăm rezultate preliminare obținute în urma procesului de adnotare și a analizei formațiilor / pseudo-formațiilor consoanelor *v* și *f* din cuvintele *vată, fată*. Graficele cu traseele F1, F2, F3, pentru toate cuvintele extrase spre analiză au fost efectuate pentru toți cei patru subiecți amintiți anterior. Un exemplu de adnotare a unui fișier ce conține consoanele siflantă *f* și semi-vocală *v* în context CV este ilustrat în figura 4. Partea superioară reprezintă semnalul în domeniul amplitudine-timp („forma de undă”), partea mediană reprezintă traseele formațiilor (pseudo-formațiilor) F1, F2, F3, F4, iar partea inferioară adnotarea la nivel de fonem (1), silabă (2) și cuvânt (3).

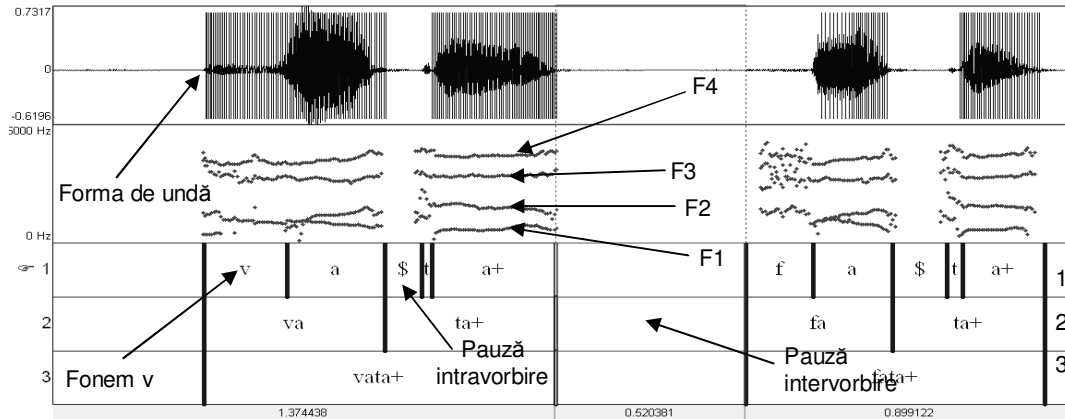


Figura 4: Exemplu de adnotare în utilitarul Praat™ (Vorbitor 1234m).

În graficul a) din figura 5 sunt reprezentate traseele formanților (pseudo-formanților) F1 și F2 ai consoanelor *v* și *f* în context CV (cuvintele *vată*, *fată*), pentru subiectul cu edentație majoră 5343m. La consoana *v*, se observă că formantul F1 are valori de cca. 250 Hz și formantul F2 de cca. 1000 Hz, la aproximativ prima jumătate a traseului, ulterior valorile acestora modificându-se brusc la 1000 Hz respectiv 2000 Hz.

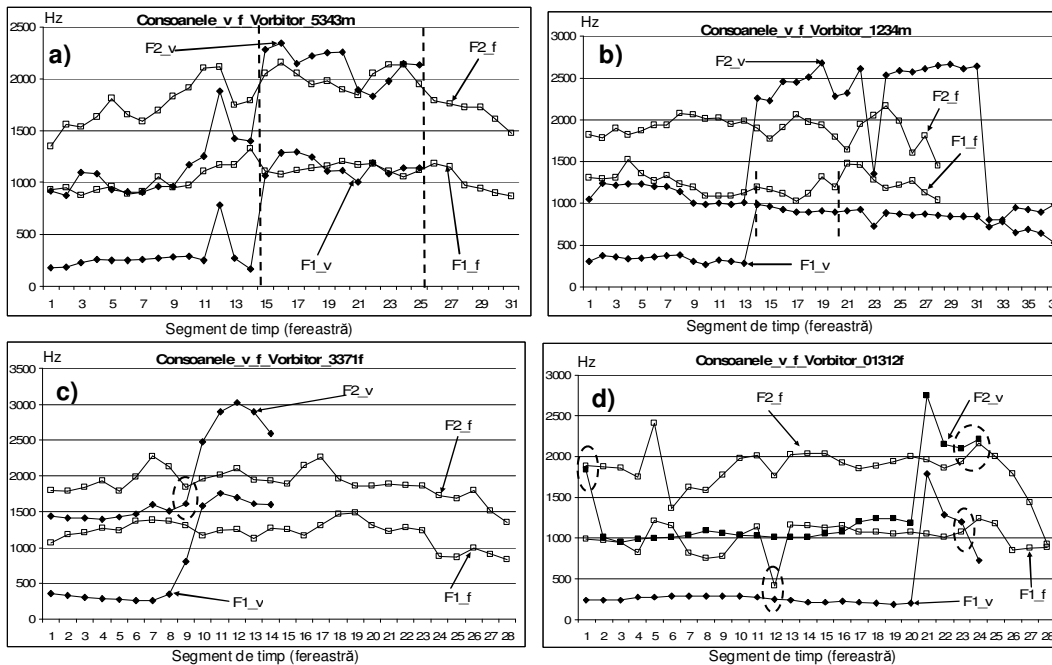


Figura 5: Traseele formanților / pseudo-formanților F1, F2, pentru consoanele *v* și *f* din cuvintele *vată*, *fată*, pentru patru vorbitori -5343m, 1234m, 3371f, 01312f ce prezintă patologiiile amintite anterior (vezi cap. „Statistica patologiilor”).

Se constată că valorile pseudo-formanților consoanei *f* au valori foarte apropiate cu cele ale formanților din a doua jumătate a traseelor consoanei *v*. Prin urmare consoana *v* alunecă prin valorile formanților către un *f*. Această alunecare poate fi interpretată ca o deficiență în pronunția consoanei *v*, dar o analiză de detaliu este necesară pentru a confirma această concluzie preliminară. În graficele b), c) și d) sunt reprezentate

traseele formanților (pseudo-formanților) consoanelor *v* și *f* pentru subiecții 1234m, 3371f și 01312f, care nu au defecte majore de dentiție. Comparativ cu subiectul ce prezintă edentație, se observă că la ceilalți trei subiecți traseele formanților F1 și F2 pentru consoanele *v* și *f* diferă între ele cu excepția zonelor marcate cu linii punctate, care cuprind: un eşantion pentru formantul F2 la subiectul 3371f, două eşantioane pentru F1 respectiv trei eşantioane pentru F2 la subiectul 01312f și șapte eşantioane pentru formantul F1 la subiectul 1234m.

6. Concluzii și direcții viitoare

Rezultatele obținute indică o posibilă corelare între modificările de dentiție și caracteristicile formantice la pronunția consoanelor fricative. Având în vedere că sunt foarte preliminare este necesară o bază statistică mai mare pentru a deriva concluzii solide. Analiza gnatofonică se află încă în stadiul de cercetare, dar în viitor poate constitui un instrument util în stomatologie. Efectuarea de înregistrări gnatofonice pe pacienți cu defecte de dentiție de tip edentație, malocluzie, disfuncții ale articulației temporo-mandibulare, înainte de aplicarea tratamentului, poate avea relevanță în stabilirea diagnosticului și evaluarea deficienței de pronunție indusă de patologiiile aparatului stomatognat. Corpusul de sunete gnatofonice poate fi utilizat și în alte aplicații medicale cum ar fi detectarea defectelor de pronunție care apar în cazul unor patologii respiratorii, laringeale, neurologice, psihologice, în măsura în care va include și astfel de patologii. O altă aplicație constă în stabilirea eficienței tratamentului logopedic, prin analiza articulării cuvintelor înainte și după aplicarea acestuia. Ca direcție viitoare ne propunem să dezvoltăm corpusul de sunete gnatofonice adnotate, prin realizarea de înregistrări atât pe subiecți cu patologii ale aparatului stomatognat, cât și pe subiecți cu dentiție normală.

Contribuția autorilor. Primul autor a conceput structura și conținutul lucrării, metoda de înregistrare și preprocesare a înregistrărilor, metodologia de analiză formantică și temporală, a dictat parte din informațiile prezentate etc. Al doilea autor a realizat înregistrările folosite aici, pe care le-a filtrat, împărțit în subfișiere și adnotat, a completat fișele subiecților, a extras valorile temporale și ale formanților, a realizat graficele cu traseele formanților pentru consoanele *v*, *f* din cuvintele analizate pentru patru vorbitori. Ambii autori au contribuit la redactarea lucrării.

Mulțumiri. Autorii sunt recunoscători persoanelor care au acceptat să fie înregistrate și au fost de acord cu utilizarea înregistrărilor în scop de cercetare și includerea acestora în baza de date existentă pe sit-ul „Sunetele limbii române”, în condițiile protecției și confidențialității datelor personale. Mulțumim întregului colectiv SROl și în mod special colegelor M. Feraru și L. Pistol pentru contribuțiile aduse la „Arhiva de sunete gnatofonice și gnatosonice” de pe sit, și colegului M. Zbancioc pentru observațiile pertinente.

Referințe bibliografice

- Boersma, P.P.G. (2002). Praat, A System for Doing Phonetics by Computer. Glot International, Vol. 5 No. 9/10, 341-345., <http://www.fon.hum.uva.nl/praat/>, data accesării: 1.04.2010.
- Burlui, V., Teodorescu, H.N., Morărașu C.S. (1993 a). L'analyse en fréquence de la fonction phonétique chez l'édenté total, 7^{me} Symposium Européen Sur Le Traitement de L'Edéntation Totale, Lyon, France, 1993 (vol. de rezumate al Simpozionului).

- Burlui, V., Teodorescu, H.N., Morărașu, C.S. (1993 b). L'analyse en fréquence de la restauration de la fonction phonétique chez l'édenté total, *Symposium Européen Sur Le Traitement de l'Édentation Totale*, Lyon, Franța.
- Burlui, V., Teodorescu, H.N., Morărașu, C.S. (1994). La fonction phonatoire chez l'édenté total. Analyse en fréquence. *Les Cahiers de Prothèse (France)*, no. 88, Decembre 1994, pp. 63-68.
- Burlui, V., Teodorescu, H.N., Morărașu, C. (1996). Analiza în frecvența gnatoaprotetică, Sesiunea jubiliară „30 de ani de învățământ stomatologic ieșean”, 1 martie 1996.
- Johnson, N.C.L., Sandy, J.R. (1999). Tooth position and speech—is there a relationship?, *The Angle Orthodontist*, vol. 69, Issue 4, 306-310, August 1999.
- Rathbone, J.S., Snidecor, J.C. (1959). Appraisal of Speech Defects In Dental Anomalies With Reference To Speech Improvement, *The Angle Orthodontist*, vol. 29, no. 1, 54-59, 1959.
- Teodorescu, H.N., Burlui, V., Leca, P.D. (1986). Gnathosonic analyzer, *Med Biol Eng Comput.* July, 1988, 26(4):428-31.
- Teodorescu, H.N., Feraru, S.M., Trandabăț, D., Zbancioc, M., Luca, R., Verbuță, A., Hnatiuc, M., Ganea, R., Voroneanu, O., Pistol, L., Șcheianu, D. (2005-2007), situl Web Sunetele Limbii Române. http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.html
- Teodorescu, H.N., Trandabăț, D., Feraru, S.M., Zbancioc, M., Luca, R. (2007). A corpus of the sounds in the Romanian spoken language for language-related education, Chapter Six, pp. 73-90. În volumul Carlos Periñán Pascual (Editor), „Revisiting Language Learning Resources”, Cambridge Scholars Publishing (CSP), UK, 2007.
- Teodorescu, H.N., Feraru, S.M. (2007 a). A study on Speech with Manifest Emotions, *10th International Conference on Text, Speech and Dialogue, TSD 2007*, Pilsen, Czech Republic, September 3-7, 2007, *Lecture Notes in Computer Science*, Springer Verlag, vol. 4629/2007, 254-262.
- Teodorescu, H.N., Feraru, S.M. (2007 b). Micro-corpus de sunete gnatosonice și gnatofonice, Pistol, Cristea, Tufiș (Eds.), *Resurse lingvistice și instrumente pentru prelucrarea limbii române*, Editura Universității „Al. I. Cuza” Iași, 21-30, 2007.
- Teodorescu, H.N., Feraru, S.M. (2008). Classification in Gnathophonics - Preliminary Results, The Second International Symposium on Electrical and Electronics Engineering, ISEEE, 12-13 septembrie, Galați, România, Galați University Press, 525-530, 2008.

DISPOZITIV ELECTRONIC DE PREPROCESARE ÎN REGIM PARALEL A SPECTRULUI VOCAL

HULEA MIRCEA¹, UNTU ALINA²

¹*Universitatea Tehnică „Gheorghe Asachi”, Facultatea de Automatică și Calculatoare
Iași – România;*

²*Universitatea Tehnică „Gheorghe Asachi”, Facultatea de Electronică,
Telecomunicații și Tehnologia Informației, Iași – România;*

mhulea@tuiasi.ro, auntu@etti.tuiasi.ro

Rezumat

Lucrarea de față prezintă un dispozitiv hardware bazat pe amplificatoare operaționale a cărui funcție principală reprezintă discriminarea frecvențelor spectrului vocal în funcție de formanții specifici vocalelor limbii române. Proiectarea acestui sistem s-a realizat în scopul generării simultane în regim paralel a trenurilor de impulsuri ce stimulează o rețea neuronală analogică utilizată în domeniul recunoașterii vocale independente de vorbitor. Astfel, prototipul prezentat în această lucrare poate fi utilizat până în acest moment pentru evidențierea canalelor de frecvență ce semnalizează recepția vocalelor. Având în vedere rezultatele obținute în acest sens, cercetările ulterioare vor avea ca scop testarea principiilor de funcționare ale acestui sistem și pentru evidențierea canalelor de frecvență activate la recepția consoanelor.

1. Introducere

Un domeniu de o importanță majoră pe plan științific reprezintă procesarea limbajului natural de către structurile de calcul artificiale. Punctul de plecare al activităților desfășurate în această direcție de cercetare constituie recunoașterea vocală a cuvintelor care se poate realiza în mod dependent sau în mod independent de vorbitor. Cea de-a doua modalitate de recunoaștere vocală constă în extragerea și apoi recunoașterea unui set de caracteristici spectrale specifice informației transmise prin vorbire. Extragerea proprietăților ce conduc la recunoașterea vocală independentă de vorbitor a cuvintelor poate fi realizată la nivel software prin elaborarea de algoritmi sau la nivel hardware prin proiectarea de circuite electronice. Scopul acestei lucrări constă în prezentarea principiilor structurale și funcționale ale unui dispozitiv hardware ce realizează preprocesarea semnalului audio în vederea recunoașterii vocalelor limbii române. Structura prototipului a fost elaborată în vederea utilizării acestuia ca element de intrare pentru rețelele neuronale analogice ce au ca unitate funcțională neuronul electronic de inspirație biologică. Funcționarea acestui sistem are la bază principiile de percepție a sunetelor evidențiate la nivelul urechii umane care realizează divizarea spectrului audio pe canale de frecvență (Stevens, 1999). În această direcție, literatura de specialitate prezintă aspecte privitoare la existența unor elemente de preprocesare care au fost proiectate ținând cont de caracteristicile limbii engleze (Hopfield, Brody, 2001), (Wills, 2004), (Loizou, Dorman, 2006). Sistemul prezentat în cadrul acestei lucrări păstrează principiul de divizare pe canale de frecvență a spectrului vocal cu deosebirea că, filtrarea semnalului audio se realizează la nivel hardware folosind un set de parametri ce

permite captarea caracteristicilor spectrale ale vocalelor limbii române. În faza inițială, cercetările proprii efectuate în acest domeniu au folosit un element de procesare a semnalului ce avea la bază filtre rezonante trece-bandă. Performanțele acestora în delimitarea canalelor de frecvență au fost reduse deoarece prezentau o atenuare scăzută a frecvențelor din afara benzii de trecere. Efectul acestui dezavantaj crește invers proporțional cu frecvența de rezonanță a filtrelor care la valori mai mici de 1 kHz necesitau creșterea substanțială a capacităților folosite. Mai mult decât atât, pentru creșterea factorului de calitate a etajelor de filtrare și păstrarea valorilor condensatoarelor în limite acceptabile, a fost necesară utilizarea de bobine speciale cu raportul L/R ridicat. De asemenea folosirea filtrelor rezonante pentru procesarea în regim paralel a semnalului audio are ca efect negativ inductanța mutuală dintre bobine.

Prin urmare, în vederea eliminării dezavantajelor induse de filtrele rezonante, pentru obținerea canalelor de frecvență sistemul actual folosește filtre active de tip Chebyshev a căror ieșire este adaptată pentru stimularea rețelelor neuronale analogice. Frecvențele de tăiere ale acestor filtre trece-bandă s-au calculat prin utilizarea rezultatelor obținute pentru formanții F1 și F2 ai vocalelor limbii române, publicate de către H.N. Teodorescu și M. Feraru pe sit-ul „Sunetele limbii române” (SRoL) prezentat în (Teodorescu et al., 2005-2007), (Teodorescu et al., 2007), (Teodorescu, Feraru, 2007). Dintre aceste date statistice s-au folosit pentru această lucrare valorile medii ale formanților pentru vocale izolate sau în context, extrase de către autori cu utilitarul PraatTM.

În continuare, secțiunea a doua a lucrării va oferi detalii privitoare la caracteristicile funcționale ale modulelor sistemului, precum și la modul de efectuare a statisticii ce a avut ca rezultat obținerea parametrilor elementelor de filtrare a semnalului vocal. Rezultatele obținute în urma testării comportamentului sistemului în cazul rostirii susținute a vocalelor sunt evidențiate în secțiunea a treia, iar concluziile referitoare la performanțele sistemului sunt prezentate în secțiunea a patra.

2. Structura sistemului

Rețelele neuronale analogice au la bază neuroni electronici ce realizează procesarea informațiilor în regim paralel folosind trenuri de impulsuri. Atât frecvența acestor impulsuri cât și energia generată de fiecare impuls pot varia în funcție de gradul de antrenare a rețelei. Pentru ca datele primite din exterior să poată fi procesate cu ajutorul structurilor neuronale de acest tip este necesară conversia în trenuri de impulsuri a semnalelor cu variație continuă. De asemenea, datorită faptului că neuronul artificial de inspirație biologică este prin esență un detector de coincidențe temporale, un interes deosebit trebuie acordat și momentelor apariției stimulilor. Astfel, rețeaua neuronală analogică va fi sensibilă atât la frecvența stimulilor de intrare cât și la gradul de concurență a acestora. Pentru stimularea unei astfel de rețele neuronale ce poate fi antrenată pentru recunoașterea vocală independentă de vorbitor este necesară o procesare prealabilă a semnalului audio. Figura 1 prezintă elementele funcționale ale sistemului de preprocesare a semnalului vocal recepționat de la un microfon. Elementul de preamplificare *PreAmp* realizează amplificarea și filtrarea semnalului audio cu scopul obținerii benzii de frecvențe corespunzătoare spectrului vocal. Semnalul obținut este normalizat de către modulul *DynAmp* prin adaptarea dinamică a amplitudinii

semnalului cu ajutorul circuitului MAX9756. Etapa următoare a procesării semnalului este îndeplinită de modulul *FTB* ce realizează divizarea pe canale de frecvență a spectrului vocal în funcție de formații specifici vocalelor limbii române.

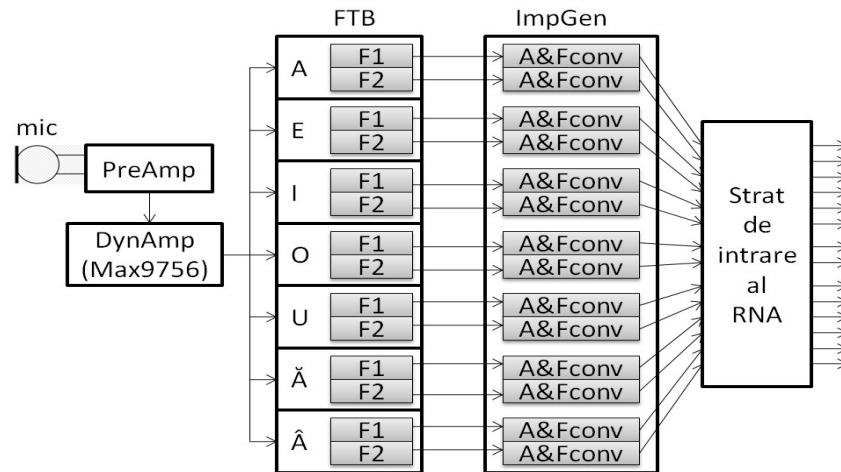


Figura 1: Structura dispozitivului de preprocesare a semnalului audio alcătuită din modulele *PreAmp* (preamplificarea semnalului), *DynAmp* (reglarea amplitudinii semnalului), *FTB* (divizarea pe canale de frecvență a spectrului vocal), *ImpGen* (conversia amplitudinii și frecvenței semnalului); Dispozitivul stimulează stratul de intrare al rețelei neuronale analogice (RNA).

Acest modul constă într-o serie de filtre active de ordin patru ce au fost proiectate folosind amplificatoare operaționale $\beta A741$. Astfel, $F1[v]$ și $F2[v]$ unde $v \in \{a, e, i, o, u, \text{ă}, \text{â}\}$ reprezintă filtre de tip Chebyshev acordate pe frecvențele specifice formațiilor $F1$ și respectiv $F2$ ai fiecărei vocale. Pentru a fi posibilă funcționarea normală a rețelei neuronale a fost necesară proiectarea unui modul auxiliar *ImpGen* care, cu ajutorul stratului de intrare a rețelei neuronale, realizează adaptarea amplitudinii și frecvenței impulsurilor generate de fiecare canal de frecvență al *FTB*.

2.1 Modulul *DynAmp*

Experimentele anterioare au evidențiat faptul că pentru obținerea unei acuități ridicate în recunoașterea vocală a fonemelor limbii române este necesară rostirea acestora cu o intensitate ridicată și aproximativ constantă. În vederea soluționării acestei probleme sistemul își propune normalizarea la nivel hardware a semnalului recepționat de la microfon după preamplificarea prealabilă a acestuia. Această adaptare de semnal se realizează cu ajutorul circuitului MAX9756 produs de firma **MAXIM**, ce îndeplinește funcția de reglare automată a amplitudinii semnalului. Prin utilizarea acestui circuit se reduce considerabil sensibilitatea răspunsului sistemului la variația intensității semnalului vocal de intrare. În continuare se vor prezenta parametrii utilizați în proiectarea schemei electronice de control a circuitului pentru inițializarea funcției de control a amplitudinii semnalului ALC (*Automatic Level Control*). Prin această funcție se amplifică componentele de intensitate scăzută ale semnalului fără a distorsiona componentele de amplitudine ridicată (MAXIM Innovation DeliveredTM, 2006). Sesizarea depășirii unui anumit prag de putere pe ieșirea circuitului determină activarea ALC ce reduce amplificarea semnalului cu până la 6 dB. Figura 2 (a) ilustrează etapele

acțiunii ALC care sunt *creșterea* în timpul t_a a atenuării puterii semnalului, *blocarea* atenuării la valoarea curentă pe parcursul intervalului de timp t_h și *revenirea* treptată a amplificării la valoarea normală în timpul t_r .

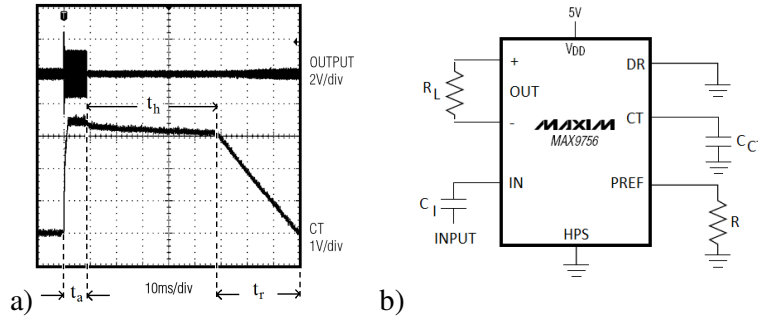


Figura 2. a) digrama de semnal obținută în urma testelor efectuate de producător asupra circuitului MAX9756 ce permite evidențierea etapelor reglării automate a amplitudinii semnalului; b) schema de conectare a pinilor prin care se definesc parametrii funcției ALC (DR, CT, PREF), precum și a pinilor ce definesc intrarea, ieșirea și selecția ieșirii (IN, OUT și HPS).

Circuitul integrat oferă posibilitatea ajustării timpului de atac t_a și a timpului de revenire t_r prin modificarea parametrilor C_{CT} și DR . Având în vedere scopul pentru care este proiectat dispozitivul, s-a ales valoarea minimă pentru $C_{CT} = 10nF$ deoarece dezavantajele vitezei ridicate de răspuns a ALC (MAXIM Innovation Delivered™, 2006) nu sunt sesizabile în cazul recunoașterii vocale. De asemenea, s-a ales $t_r = 300ms$ (maxim) prin conectarea DR la GND pentru a preveni fluctuațiile rapide de intensitate a semnalului vocal pe parcursul cuvintelor.

2.2 Modulul FTB

Pentru a fi posibilă utilizarea rețelelor neuronale analogice în domeniul procesării artificiale a semnalului vocal este necesar ca receptorii specializați în detecția caracteristicilor spectrale ale fonemelor să aibă capacitatea de a fi activați simultan. Acest fapt prezintă o importanță majoră datorită faptului că în mod normal un fonem este caracterizat de un număr de frecvențe a căror prezență în spectrul vocal trebuie semnalizată rețelei neuronale în timp real. Recunoașterea fonemului respectiv se produce cu condiția ca apariția acestor frecvențe la intrarea rețelei neuronale să se producă în mod simultan. Plecând de la această idee cât și de la modelul urechii umane, modulul FTB realizează divizarea pe canale de frecvență a spectrului audio. Scopul acestei abordări reprezintă extragerea în timp real a frecvențelor ce conduc la recunoașterea informației transmise prin vorbire în mod independent de caracteristicile vocii.

2.2.1 Simularea soft a filtrelor Chebyshev

În vederea implementării etajelor de filtrare utilizate în recunoașterea vocală independentă de vorbitor, s-au proiectat cu ajutorul utilitarului *FilterLab* oferit gratuit de *Microchip*, câte două filtre trece-bandă de tip Chebyshev pentru fiecare vocală a limbii române. Având în vedere faptul că formanții F1 și F2 (multipli ai frecvenței

fundamentale) sunt implicați în discriminarea vocalelor, pentru recunoașterea acestora s-a implementat câte un filtru ce evidențiază recepția fiecăruia dintre cei doi formanți. Parametrii de proiectare utilizați sunt: tip filtru *Chebyshev*; riplu permis în banda de trecere -3 dB; numărul de poli 4; câștig 1. De asemenea frecvențele de tăiere pentru cele șapte vocale (*a, e, i, o, u, ă, â*) sunt prezentate în tabelul 1.

2.2.2 Obținerea frecvențelor de tăiere

Pe baza valorilor medii ale formanților preluate de pe sit-ul SRoL, s-a calculat media și deviația standard pentru formanții F1 și F2 specifici celor șapte vocale. Pentru stabilirea lățimii benzii de frecvență a celor două filtre corespunzătoare fiecărei vocale, am ales drept frecvență de tăiere minimă f_{\min} și maximă f_{\max} după cum urmează:

$$f_{\min} = f_{\text{central}} - st_dev \quad (1)$$

$$f_{\max} = f_{\text{central}} + st_dev \quad (2)$$

unde f_{central} este frecvența centrală a benzii de trecere, iar st_dev este deviația standard. Atât frecvența centrală cât și deviația standard s-au obținut în urma statisticii efectuate pe un număr de persoane aflate în diverse stări emoționale.

Analiza statistică s-a efectuat în mod diferit pentru grupurile de vocale *a, e, i* și respectiv, *o, u, ă, â* deoarece am avut disponibile valori ale formanților specifice vocalelor rostite în contexte diferite. Pentru *a, e, i* au existat valori pentru vocale izolate, iar pentru *o, u, ă, â* au existat valori pentru vocale în context. Cele din urmă au fost extrase din propozițiile „*Aseară*” și „*Vine mama*” pentru patru stări emoționale: neutru, fericire, furie, tristețe. Frecvența centrală a filtrelor s-a obținut prin calculul mediilor a 30 de valori pentru vocalele izolate *a, e, i*, iar pentru vocalele în context *o, u, ă, â*, media a fost calculată folosind 40 de valori. La ultimul grup de vocale s-au considerat 10 vorbitori (cinci de gen feminin și cinci de gen masculin) pe patru stări emoționale. Rezultatele obținute sunt ilustrate în tabelul 1.

Tabel 1: Valorile frecvențelor medii, minime și maxime corespunzătoare formanților vocalelor

Vocală	a		e		i		o		u		ă		â	
	F1	F2	F1	F2	F1	F2	F1	F2	F1	F2	F1	F2	F1	F2
f_{central}	751	1274	583	1889	399	2328	466	915	403	788	542	1433	389	1648
f_{\min}	618	1093	486	1671	338	2091	361	794	333	708	432	1192	325	1428
f_{\max}	883	1456	680	2107	460	2565	570	1036	473	868	651	1674	453	1867

Schema generală a unui filtru trece bandă de tip Chebyshev alcătuit din două amplificatoare operaționale este dată în figura 3.

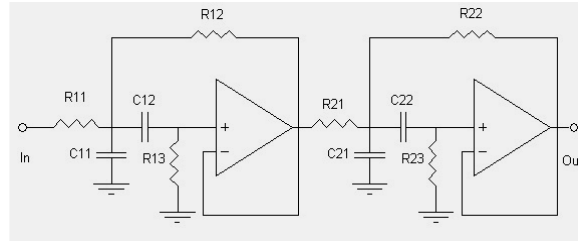


Figura 3: Schema generală a unui filtru trece bandă Chebyshev.

Pentru evidențierea comportamentului canalelor de frecvență implementate, în figura 4 este prezentată simularea în *FilterLab* a răspunsului în frecvență a filtrului trece-bandă pentru formantul F2 al vocalei *e*. Frecvențele de tăiere corespunzătoare la -3dB sunt $f_{\min} = 1671\text{Hz}$ și $f_{\max} = 2107\text{Hz}$, iar lățimea benzii de trecere este egală cu dublul deviației standard

$$\Delta f_b = 2 \cdot st_dev = 436\text{Hz} \quad (3)$$

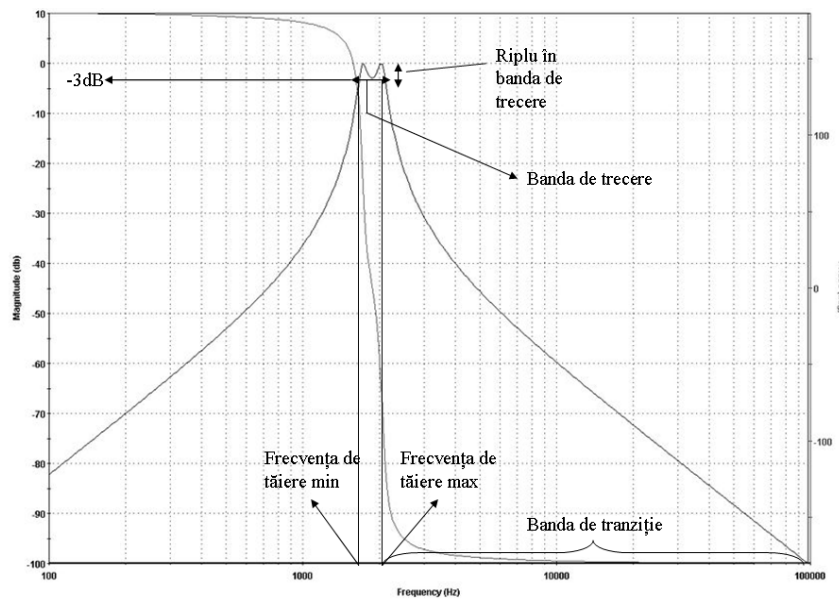


Figura 4: Răspunsul în frecvență a filtrului trece-bandă de tip Chebyshev.

Filtrele Chebyshev de tipul I sunt filtre polinomiale (numai cu poli, fără zerouri), având o caracteristică de modul cu ripluri egale în banda de trecere și monoton descrescătoare în banda de oprire. Dintre toate filtrele polinomiale de ordinul N , filtrele Chebyshev de tipul I au zona de tranziție cea mai îngustă. Performanțele filtrului sunt complet precizate de mărimea riplului în banda de trecere și de ordinul filtrului ce determină panta caracteristicii de modul. Utilizând programul *FilterLab* se pot proiecta trei tipuri de filtre: Bessel, Butterworth și Chebyshev. Filtrul Bessel se poate realiza doar pentru tipul de răspuns trece-jos, iar celelalte două pentru răspunsuri de tip trece-jos, trece-sus, și trece-bandă. În comparație cu filtrul Butterworth, filtrul Chebyshev realizează o mai bună atenuare a frecvențelor din afara benzii de trecere, acest fapt justificând alegerea făcută în procesul de proiectare a filtrelor.

2.2.3 Delimitarea canalelor de frecvență

Semnalul sinusoidal generat de filtrele trece-bandă prezentate anterior este convertit în impulsuri prin utilizarea unui comparator LM324 pentru fiecare formant. Acest etaj permite și reglarea amplitudinii de la care activitatea canalului de frecvență va produce stimularea rețelei neuronale. Prin urmare, ieșirea generată de modulul *FTB* se materializează într-un tren de impulsuri de amplitudine, durată și frecvență variabile ca în figura 5(a). Energia acestor impulsuri este convertită în trenuri de impulsuri de amplitudine și frecvență constantă de către modulul *ImpGen* ce constituie o interfață între modulul de filtrare și stratul de intrare al rețelei neuronale.

2.3 Modulul *ImpGen*

Experimentele realizate pe parcursul cercetărilor efectuate în domeniul recunoașterii vocale au arătat că pentru creșterea stabilității în funcționare a rețelei neuronale analogice este necesară o adaptare a amplitudinii și frecvenței impulsurilor ce stimulează intrările acesteia (Hulea, 2009). Schema electrică a acestui modul conectată cu un neuron electronic specific stratului de intrare al rețelei este prezentată în figura 5 (b).

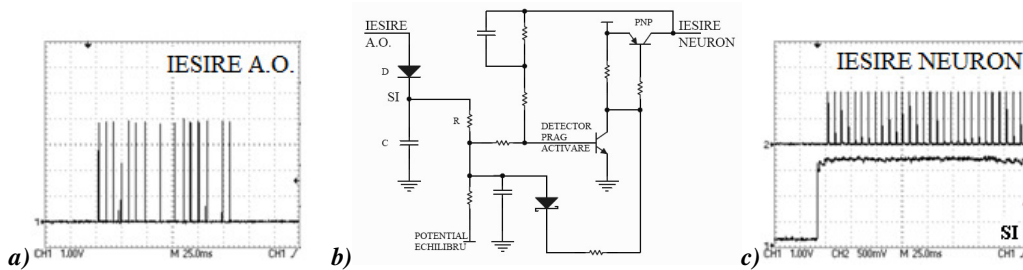


Figura 5. a) semnalul de ieșire a unui canal de frecvență al modulului *FTB*; b) schema electrică a ansamblului [modul *ImpGen* - neuron de intrare] pentru un canal de frecvență oarecare; c) ieșirea neuronului de intrare (semnalul IESIRE NEURON) și integrarea ieșirilor modulului *FTB* de către capacitatea *C* (Semnalul SI).

Conversia ieșirii modulului *FTB* constă în integrarea stimulilor generați și transformarea energiei obținute în impulsuri de amplitudine și frecvență constante. Prima etapă se realizează folosind capacitatea *C* ce se încarcă prin dioda *D*, iar cea de-a doua este asigurată de componentele electronice ce compun neuronul de intrare (Hulea, 2008). Frecvența impulsurilor de ieșire este invers proporțională cu valoarea rezistenței *R* prin care are loc descărcarea condensatorului *C*. În cazul sistemului descris în această lucrare, valoarea de 470 k Ω a acestei rezistențe a fost aleasă pentru a obține la ieșire o frecvență de aproximativ 160 Hz.

3. Rezultate experimentale

Răspunsurile în frecvență ale filtrelor de tip Chebyshev obținute în urma simulărilor efectuate în programul *FilterLab* au arătat că performanța acestui tip de filtru se

îmbunătățește proporțional cu ordinul acestuia. Prin urmare, odată cu creșterea numărului de etaje de filtrare utilizate, crește și panta caracteristicii de modul a filtrului. Pe de altă parte, panta maximă de tranziție se obține prin delimitarea strictă a spectrului util cu ajutorul comparatoarelor LM324 implementate pe ieșirea fiecărui canal de frecvență. Stabilitatea limitelor benzii de trecere obținute în acest mod este asigurată în cazul în care variația amplitudinii semnalului este redusă la minim. Această condiție este îndeplinită prin reglarea dinamică a amplitudinii semnalului de către modulul *DynAmp* conectat pe intrarea modului *FTB*. Pentru ilustrarea comportamentului descris anterior, figura 6 prezintă câteva aspecte legate de funcționarea canalului de frecvență centrală 542 Hz corespunzător formantului F1 al vocalei /ă/. Formele de undă prezentate au fost obținute în urma evidențierii cu ajutorul osciloscopului a semnalelor generate prin nodurile de conectare a celor patru module ale sistemului. Vocalele au fost rostite independent de context la o distanță de aproximativ 3cm de microfon.

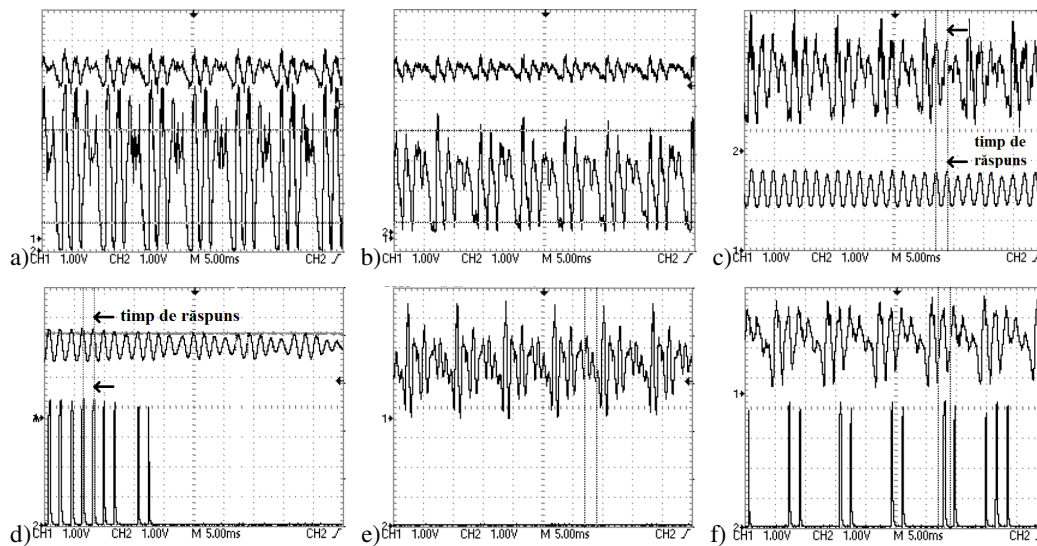


Figura 6. a) Intrarea (sus) și ieșirea (jos) a modului *DynAmp* când funcția de reglare automată a amplitudinii semnalului (ALC) este dezactivată; b) intrarea (sus) și ieșirea (jos) a modului *DynAmp* după activarea ALC; c) ieșirea *DynAmp* (sus) și ieșirea filtrului Chebyshev corespunzător formantului F1 (jos) în timpul rostirii vocalei susținute a vocalei /ă/ (cursorii verticali evidențiază timpul de răspuns a *FTB*); d) Semnalul generat de filtrul Chebyshev (sus) și ieșirea modului *FTB* corespunzătoare formantului F1 (jos) a vocalei /ă/ la rostirea vocalei /ă/; semnalul continuu reprezintă pragul de la care semnalul generat de canalul de frecvență este considerat util (cursorii verticali evidențiază timpul de răspuns a *ImpGen*); e) ieșirea circuitului MAX9756 la rostirea vocalei /ă/ și ieșirea *FTB* corespunzătoare formantului F1 al vocalei /ă/ f) semnalul generat de *DynAmp* la rostirea vocalei /ă/ și ieșirea *FTB* corespunzătoare formantului F1 al vocalei /ă/.

Așadar primele două diagrame de semnal (a) și (b) din figura 6, evidențiază importanța funcției ALC a circuitului MAX9756 în adaptarea amplitudinii semnalului generat de modulul *PreAmp*. Semnalele prezentate în figura 6 (c) prezintă efectul produs de filtrul Chebyshev de frecvență centrală 542 Hz corespunzător formantului F1 al vocalei /ă/ în momentul rostirii susținute a acesteia. Diagrama de semnal (d) ale aceleiași figuri prezintă pragul de la care activitatea canalului de frecvență va fi considerată utilă. În

continuare sunt prezentate rezultatele unui test efectuat asupra elementelor de filtrare corespunzătoare canalului de frecvență menționat. Pentru efectuarea acestui test s-a rostit vocala /a/ și s-a verificat activitatea ieșirii *FTB* corespunzătoare formantului F1 al vocalei /ă/. După cum se poate observa, rezultatele testului au fost cele așteptate în sensul că vocala /a/ nu a produs activarea canalului de 542 Hz (diagrama de semnal (e)), acesta fiind activat la recepția vocalei /ă/ (diagrama de semnal (f)).

Valorile obținute pentru frecvențelor de tăiere ale filtrelor (tabelul 1) au impus efectuarea unor ajustări în implementarea dispozitivului pentru eliminarea elementelor de filtrare redundante. În acest sens, pentru formantul F1 al vocalelor *i, u*, și *â* s-a utilizat același filtru de frecvență centrală 400 Hz, iar pentru formantul F1 al vocalelor *e* și *ă* s-a utilizat un filtru de frecvență centrală 560 Hz. Diferența de aproximativ 200 Hz dintre valorile medii obținute pentru formantul F2 fac posibilă discriminarea vocalelor limbii române folosind doar filtrele corespunzătoare acestui formant.

Un alt aspect important ce a fost luat în considerare în procesul de testare a sistemului reprezintă evaluarea timpului de răspuns a acestuia. Cursorii verticali de pe diagramele de semnal (c) și (d) ale figurii 6 evidențiază timpii de răspuns a submodulelor *FTB* și respectiv *ImpGen* care sunt neglijabili în comparație cu perioada semnalelor specifice spectrului vocal. De asemenea diagramele de semnal (a) și (b) ale aceleiași figuri evidențiază cu ușurință timpul de răspuns a modulului *DynAmp* care de asemenea este neglijabil. Prin urmare se poate deduce faptul că timpul total de răspuns a întregului sistem de preprocesare este neglijabil în comparație cu variația semnalelor specifice spectrului vocal.

4. Concluzii

Pentru a beneficia de avantajele utilizării rețelelor neuronale cu răspuns în timp real în domeniul recunoașterii vocale, s-a proiectat un sistem de preprocesare în regim paralel a semnalului audio. Plecând de la modelul aparatului auditiv uman ce realizează divizarea spectrului vocal pe canale de frecvență s-a implementat un set de filtre trece-bandă acordate pe frecvențele caracteristice formațiilor vocalelor limbii române. În vederea obținerii unei acuități ridicate în recunoașterea vocală independentă de vorbitor s-a urmărit creșterea performanțelor sistemului de filtrare. Aceasta s-a realizat atât prin utilizarea filtrelor de tip Chebyshev care prezintă cea mai bună atenuare a frecvențelor din afara benzii de trecere, cât și prin folosirea circuitului MAX9756 care realizează adaptarea dinamică a amplitudinii semnalului. De asemenea, procesul de selecție a frecvențelor de tăiere pentru aceste filtre s-a bazat pe date statistice ce cuprind o serie de voci de gen masculin și feminin precum și diferite stări emoționale.

Experimentele efectuate pentru studiul funcționării sistemului prezentat în această lucrare au arătat o superioritate evidentă a performanțelor filtrelor Chebyshev în comparație cu filtrele rezonante utilizate în cadrul cercetărilor anterioare descrise în (Hulea, 2008). Performanțele sistemului s-au îmbunătățit substanțial mai ales în cazul canalelor de frecvențe centrale mai mici de 1000 Hz. Un alt avantaj al filtrelor Chebyshev reprezintă creșterea eficienței filtrelor prin introducerea elementelor active (A.O.) și eliminarea componentelor inductive (L). De asemenea capacitățile folosite s-au redus cu trei ordine de mărime ceea ce reduce costurile implementării. De asemenea, îmbunătățirea performanțelor în filtrare a fost obținută fără a afecta răspunsul în timp

real a sistemului. Totuși o consecință a utilizării circuitului MAX9756 pentru adaptarea amplitudinii semnalului audio reprezintă creșterea consumului de curent al dispozitivului cu aproximativ 48 mA. Cu toate acestea, acest aspect trebuie luat în considerare doar în cazul utilizării unei surse portabile de energie pentru alimentarea sistemului. Având în vedere creșterea acuității în delimitarea benzilor de frecvență a sistemului de filtrare se preconizează pe de o parte îmbunătățirea performanțelor în domeniul recunoașterii vocale independente de vorbitor, cât și simplificarea prin reducerea numărului de neuroni a arhitecturii rețelei neuronale utilizate în detecția cuvintelor.

Referințe bibliografice

- Hopfield, J. J., Brody C.D. (2001). What is a moment? Transient synchrony as a collective mechanism for spatiotemporal integration, *Proceedings of the National Academy of Sciences of the USA*, 1282 – 1287.
- Hulea, M. (2008). A Model of Silicon Neurons Suitable for Speech Recognition, *Computer Engineering and Applied Informatics*, Vol.10, Nr. 4, 32-41, ISSN 1454-8658.
- Hulea, M. (2009). A New Method to Obtain Non-Volatile Memory for Networks of Spiking Neurons, *Memoirs of the Scientific Sections, Romanian Academy* în curs de apariție.
- Loizou P.C., Dorman, M. (2006). On the Number of Channels Needed to Understand Speech, *Journal of the Acoustical Society in America*.
- MAXIM Innovation Delivered™ (2006). 2.3W Stereo Speaker Amplifiers and DirectDrive Headphone Amplifiers with Automatic Level Control.
- Stevens, K. (1999). *Acoustic Phonetics*, Editura Cambridge MA: The MIT Press, 204 - 214.
- Teodorescu, H.N., Feraru, S.M., Trandabăț, D., Zbancioc, M., Luca, R., Verbuță, A., Hnatiuc, M., Ganea, R., Voroneanu, O., Pistol, L., Șcheianu, D. (2005-2007), situl Web Sunetele Limbii Române.
http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.html
- Teodorescu, H.N., Trandabăț, D., Feraru, S.M., Zbancioc, M., Luca, R. (2007). A corpus of the sounds in the Romanian spoken language for language-related education, Chapter Six, pp. 73-90. În volumul Carlos Periñán Pascual (Editor), „Revisiting Language Learning Resources”, Cambridge Scholars Publishing (CSP), UK, ISBN 1-84718-156-2; ISBN 13: 9781847181565, 2007.
- Teodorescu, H.N., Feraru, S.M. (2007). A study on Speech with Manifest Emotions, *10th International Conference on Text, Speech and Dialogue, TSD 2007*, Pilsen, Czech Republic, September 3-7, 2007, *Lecture Notes in Computer Science*, Springer Verlag, vol. 4629/2007, 254-262, ISBN 978-3-540-74627-0.
- Wills, S. (2004). Computation with spiking neurons, *teză de dizertație trimisă la University of Cambridge pentru admiterea la doctorat*, 5-12.

CAPITOLUL 2

PLATFORME, DICȚIONARE ȘI CORPUSURI PENTRU PRELUCRAREA TEXTELOR

RESURSE LINGVISTICE ROMÂNEȘTI ÎN FLUX CONTINUU

DAN CRISTEA

Universitatea „Alexandru Ioan Cuza” din Iași

Str. Berthelot nr. 16, 700483 – Iași, Romania

dcristea@info.uaic.ro

Rezumat

Lucrarea prefigurează o inițiativă legislativă pentru obținerea de resurse lingvistice românești la scară mare, fără atingerea intereselor de proprietate intelectuală ori financiare ale autorilor și deținătorilor de texte. Conservarea limbii române scrise în vederea urmării continue a evoluției ei este perfect realizabilă din punct de vedere tehnologic. Se descrie un flux tehnologic ce are ca obiectiv colectarea tuturor textelor românești ce urmează a fi tipărite, stocarea lor pe durată indefinită și prelucrarea lor în scop științific. Sunt imaginate servicii în folosul cercetătorilor limbii cât și a proprietarilor drepturilor de autor și a autorilor de texte.

1. Introducere

Lucrarea aduce în atenție o propunere pentru colectarea, conservarea și utilizarea în scopul cercetării, a tuturor textelor publicate în România de către edituri.

Este evident faptul că, fără un efort continuu, acele limbi care, sunt cunoscute acum a avea „mai puține resurse” vor continua să fie văzute astfel chiar și atunci când, ipotetic vorbind, ar ajunge la un nivel de reprezentare în resurse echivalent nivelului actual al așa-ziselor „limbi mari”. Mai mult decât atât, dacă limbile azi bogat reprezentate în resurse vor înceta acum să achiziționeze resurse, în baza raționamentului că necesitățile lor de cercetare sunt satisfăcute, în scurt timp ele își vor pierde acest avantaj. Acest lucru se datorează unui proces de deteriorare (îmbătrânire) rapidă a resursele lingvistice. Îmbătrânirea se datorează evoluției continue a limbilor dar și modificării viziunii cercetătorilor cu privire la oglindirea fenomenelor de limbă. În cazul resurselor adnotate, faptele lingvistice, care fac obiectul adnotării automate, se pot schimba de-a lungul timpului, pe măsură ce teoriile lingvistice, pe care se bazează convențiile de marcare, evoluează și pe măsură ce însăși procesele de adnotare automată se îmbunătățesc. Putem spune că nu există un final în construirea de resurse lingvistice.

În multe țări există o lege a „depozitului legal”. Aceasta obligă toți furnizorii de materiale tipărite (edituri, persoane fizice sau juridice, care imprimă documente pentru public, casele de înregistrări și studiourile, Banca Națională, Monetăria Statului, Poșta Națională etc.) – să îi denumim *furnizori de resurse* – să trimită unei biblioteci naționale (care poate fi o singură unitate fizică sau un consorțiu de biblioteci), pentru conservare pe termen lung, un număr de copii pentru fiecare element imprimat destinat distribuției. Cu toate ca orizontul producțiilor de texte s-a schimbat dramatic în ultimii ani, formatul electronic surclasând formatul tipărit, după știința mea există doar încercări timide de actualizare corespunzătoare și a aspectelor juridice.

Pe măsură ce resursele lingvistice devin tot mai necesare pentru studiul limbii și cum obținerea lor este adesea costisitoare, problematica achiziționării lor ar trebui să nu mai fie accidentală sau episodică, ci ar trebui să devină o politică națională. Ceva trebuie făcut. Ar trebui să existe o lege care să protejeze resursele lingvistice ale limbilor ce se vorbesc într-o țară și aceasta ar trebui să fie de interes primordial. Această lucrare discută o posibilă soluție, care, cu toate că nu este ușor de implementat, ar putea să schimbe complet scena resurselor lingvistice în viitorul apropiat.

2. Îmbunătățirea legislației cu privire la depozitul legal

O investigație recentă în rândul câtorva dintre cei mai importanți producători de informație tipărită din România a scos în evidență faptul ca multe edituri ar fi dispuse să-și doneze resursele în scopul cercetării. Cu toate acestea, un alt segment, care din nefericire este majoritar, nu este interesat într-o colaborare cu oamenii de știință. Acești producători se tem că cedarea datelor lor textuale în folosul cercetării ar echivala cu pierderea controlului asupra proprietății lor, ar declanșa o diminuare a profitului, sau, pur și simplu, ignoră importanța problemei și nu au timp să se dedice acestor aspecte.

În realitate, nimic din toate acestea nu ar trebuie să se întâmple. Cu toate că avem nevoie de datele lor lingvistice, nu dorim ca *furnizorii de resurse* să fie păgubiți dacă oferă textele în beneficiul științei. Ideea este de a promova o inițiativă legislativă (să o numim **Legea depozitului electronic**) care să impună obligativitatea ca *furnizorii de resurse* să doneze datele lor lingvistice pentru supravegherea permanentă a evoluției limbii.

Următoarele tipuri de resurse, produse în serie, ar trebui să facă obiectul unei astfel de inițiative legislative, indiferent dacă resursele respective sunt destinate distribuției comerciale sau gratuite: cărți, broșuri, pliante, jurnale, reviste, almanahuri, calendare, partituri muzicale, materiale propagandistice cu scop politic, administrativ, cultural, artistic, științific, educațional, religios; postere, declarații și alte materiale destinate publicării în spații publice, teze de doctorat, cursuri universitare, documente în format electronic care conțin material lingvistic (CD-uri, DVD-uri etc.), standarde și norme tehnice, publicații emise de autorități naționale și locale, colecții de norme și legi și orice alt material imprimat sau multiplicat folosind metode grafice sau fizico-chimice.

La nivel practic, inițiativa presupune existența unui depozit central, care este o entitate (centru de calcul, institut etc. – să-l denumim generic *Portal*), care, pe de o parte, are autoritatea legală de a primi și stoca date oferite prin contribuția *furnizorilor de resurse* și, pe de altă parte, este echipat tehnic pentru a colecta, stoca pentru o perioadă de timp nelimitată și prelucra, toate textele în format electronic ce urmează a fi tipărite și distribuite, zilnic, într-o țară.

Legea ar trebui să stipuleze că, prin trimiterea unei copii electronice acestui depozit național pentru conservarea pe termen nelimitat, nici un drept de autor sau beneficiu comercial nu va fi pierdut de către *furnizorul de resurse* sau autor. Aceasta copie va putea fi folosită, prin intermediul *Portalului*, **doar în scopul cercetării lingvistice** și nici *Portalul* nici vreo altă persoană ori entitate nu vor putea face publice datele sau fragmente care să exemplifice o cercetare, în Internet sau într-un alt mediu, decât cu acordul proprietarului. Este clar că această lege va trebui să includă un capitol destinat

drepturilor de autor, care să fie extrem de atent conturat (COM, 2009). Alternativ, e posibil ca legea drepturilor de autor și a proprietății intelectuale să necesite amendamente care să o armonizeze cu conținutul noii legi a depozitului electronic. De asemenea, o slabă caracterizeze a măsurilor de securitate în transferul de date nu poate fi acceptabilă în formularea acestei legi.

În multe țări se constată un interes ridicat în privința organizării unor servicii publice de depozitare a conținutului electronic în biblioteci digitale¹. Multe dintre acestea sunt arondate instituțiilor de cercetare, scopul lor fiind cu precădere acela de a păstra conținutul științific în vederea diseminării lui cu, ori fără, restricții. Sunt extrem de numeroase apoi inițiativele de păstrare a cărților rezultate din scanări în format electronic (Galica, în Franța², proiectul Gutenberg³ cu multitudinea de proiecte afiliate, serviciul de căutare oferit de Google Books⁴, proiectul Runeberg⁵ pentru literatură nordică, consorțiul Open Content Alliance⁶ etc.⁷).

3. Fluxul colectării datelor

Cred că o metaforă adecvată pentru *Portal* ar fi: o fabrică care procesează continuu cuvinte. Activitatea *Portalului* trebuie să înceapă cu crearea de identități în Portal tuturor *furnizorilor de resurse*: înscrierea la Portal a unei entități ca *furnizor de resurse* trebuie să se efectueze la momentul intrării în acțiune a Legii sau înainte de distribuirea primei publicații. În urma înscrierii, *furnizorul de resurse* va obține un cod de identificare (*FID*) de la *Portal*. Acest cod va fi folosit pentru comunicarea cu *Portalul*, cu privire la orice publicație, pe toată perioada de activitate a acestuia.

Să presupunem că astăzi, *furnizorul de resurse* pregătește pentru publicare un nou document *D*, care a primit deja aprobarea editorială „bun de tipar”. Figura 1.a prezintă fluxul de date inițiat de acest nou element. *Furnizorul de resurse* completează un formular electronic (antet, *header* – *H*), care conține informații de identificare a documentului, și apoi interacționează cu *Portalul*, încercând *FID*-ul său, antetul *H* și o copie editabilă a documentului *D*. *Portalul* primește aceste informații și solicită un cod de identificare permanent (*persistent ID* – *PID*) de la o autoritate care îl poate elibera (Kunze and Rogers, 2003; Schwardmann, 2009). Atunci când primește un astfel de cod, stochează în depozit un pachet conținând *PID*, *FID*, *H*, și *D*. După aceasta, *Portalul* va întoarce *furnizorului de resurse* un mesaj de confirmare care conține două părți: o parte ce poate fi citită în clar și o parte ce conține un cod de bare. Căsuța de confirmare trebuie să înregistreze un sigiliu al *Portalului*, împreună cu *PID*, *FID* și *H*. Documentul, care conține acum și căsuța de confirmare din partea *Portalului*, pe o manșetă, o copertă interioară sau într-o anexă, poate fi tipărit (Figura 1.b). Această căsuță ar trebui să dovedească oricărei autorități însărcinate cu verificarea respectării Legii ca depozitul

¹ http://en.wikipedia.org/wiki/Digital_library

² <http://gallica.bnf.fr/?lang=FR>

³ http://www.gutenberg.org/wiki/Main_Page

⁴ <http://books.google.com/>

⁵ <http://runeberg.org/>

⁶ <http://www.opencontentalliance.org/>

⁷ O listă a proiectelor de biblioteci digitale poate fi consultată aici:

http://en.wikipedia.org/wiki/List_of_digital_library_projects

electronic legal al acelu document a fost efectuat, toate informațiile necesare pentru identificare fiind acolo.

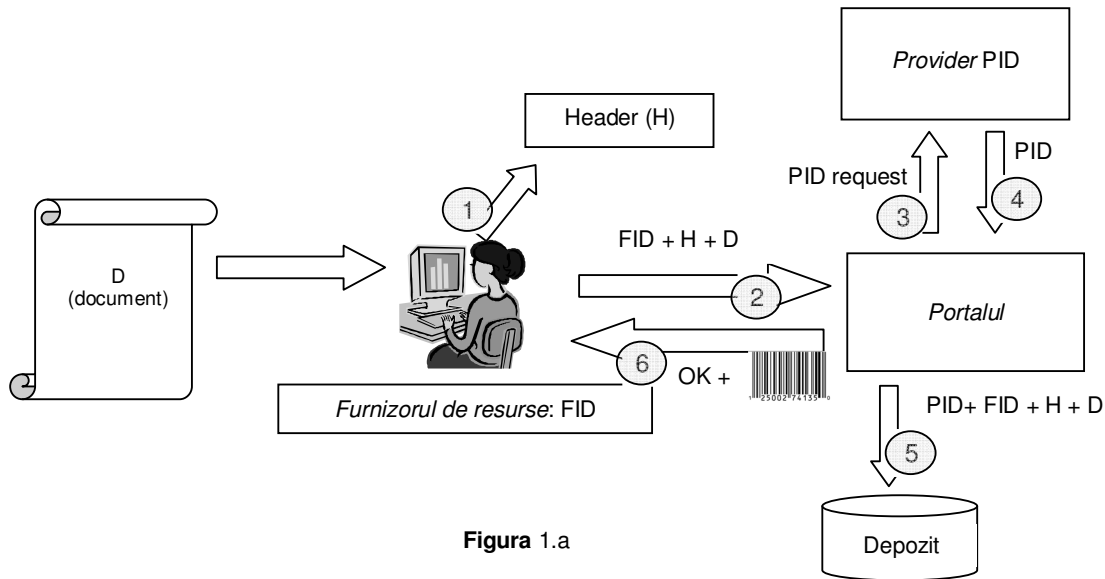


Figura 1.a

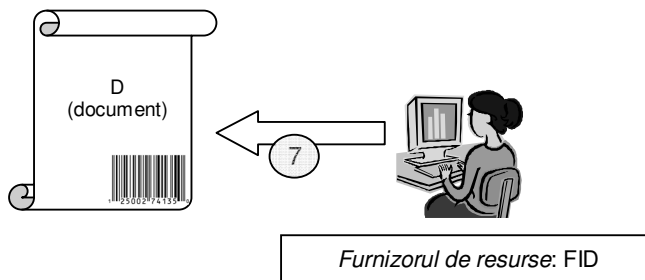


Figura 1.b

Schimbul de date detaliat mai sus între *furnizorul de resurse* și *Portal*, incluzând și o comunicare cu un terț, responsabil pentru emiterea *PID*-urilor, pare greoi și consumator de timp și, din această cauză, total neacceptabil pentru edituri. Într-adevăr, este cunoscut faptul ca activitatea în aceste unități este extrem de accelerată, ele fiind constrânse să proceseze informațiile cu mare rapiditate, adesea în condiții de stres, cu precădere atunci când tipăresc cotidiene. În fapt însă, tot fluxul de informații descris mai sus poate fi automatizat, inclusiv completarea antetului *H*, care poate fi lăsată în grija unor module specializate capabile să preia elemente de identificare din documentul electronic *D*. Practic, întregul lanț poate fi realizat la apăsarea de către editor a unui buton aflat în aplicația de editare. Rezultatul, materializat instantaneu, va fi caseta de confirmare, pe care editorul trebuie numai să o plaseze acolo unde dorește pentru a fi tipărită împreună cu documentul.

4. Fluxul de prelucrare a datelor

Odată obținute, datele din *Portal* trebuie prelucrate. În această secțiune voi descrie o listă a capacităților de procesare pe care *Portalul* ar trebui să le ofere.

Mai întâi, este evident că *Portalul* trebuie să posede capacități de stocare de foarte mari dimensiuni și că ele trebuie special proiectate pentru conservarea datelor pe perioade de timp oricât de îndelungate. Mai apoi, el trebuie să ofere capacități de indexare, căutare și acces, la diverse nivele: antet, semne lexicale (cuvinte), expresii lexicale, precum și coocurențe contextuale. Printre altele, acest lucru presupune ca fiecare document ce este plasat în *Portal* să fie supus unui lanț de procesări care trebuie să includă minimal: segmentare la nivel de cuvânt, etichetare la parte de vorbire, lematizare și indexare. Înregistrarea textului în forma inițială este de asemenea de prevăzut, pentru că la el ar trebui să facă trimitere adnotările XML *stand-off*. Este de așteptat ca adnotările XML și indecșii să multiplice dimensiunea inițială a documentelor text de câteva ori.

Pornind de la aceste funcționalități de bază, o linie diferită de prelucrări se va adresa nevoilor lexicografice. *Portalul*, funcționând, din acest punct de vedere, în același timp ca un depozit imens de resurse dar și ca o colecție de instrumente de prelucrare a limbii, va trebui să fie capabil să realizeze operațiuni complexe cum ar fi: identificarea cuvintelor străine, semnalizarea unor cuvinte noi, recunoașterea sensului cuvintelor în context, identificarea sensurilor noi, semnalizarea cuvintelor uitate (învechite), semnalizarea sensurilor care nu mai sunt folosite etc. De exemplu, se poate imagina că semnalizarea intrării în limbă a unor cuvinte noi, cât și ieșirea unor cuvinte din circulație, ar trebui declanșată de existența, permanent actualizată, a frecvenței de apariție a cuvintelor pe parcursul unui interval constant de timp, atunci când această frecvență se plasează deasupra (ori dedesubtul) anumitor praguri, decise de o autoritate lingvistică. În mod similar, semnalizarea unui nou sens poate fi declanșată de eșecul de a cataloga sensul, desprins din context, între sensurile păstrate într-un repozitoriu de sensuri, cum ar fi, de exemplu, un dicționar explicativ, dacă acest lucru are loc recent cu o anumită frecvență și dacă modelul este suficient de stabil. Similar, sensurile uitate (învechite) sunt recunoscute atunci când regăsirea acestora scade sub un anumit prag.

Procesul care trebuie pus la baza recunoașterii cuvintelor ori sensurilor învechite presupune monitorizarea permanentă a unui set de cuvinte luate în „colimator”. Aceste cuvinte/sensuri sunt plauzibile a fi candidați la dispariție deoarece apar cu frecvențe care scad continuu. Trebuie notat faptul că, de-a lungul unui interval de timp, criteriul frecvenței absolute, ori chiar al uneia relative, poate să nu fie relevant, deoarece unele cuvinte sunt foarte rar folosite, cu toate că nu sunt în pericol de a fi considerate moarte (de exemplu unele neologisme științifice). Cea mai bună cale ar fi crearea câte unei fișe „personale” pentru fiecare lexical, în care să se înregistreze un set de proprietăți (modificabile în timp), printre care și frecvența aparițiilor de-a lungul timpului (un grafic, pe baza căruia să poată fi calculat un grad de deteriorare a utilizării), lista registrelor unde se remarcă folosirea lui (cu frecvențele relative asociate) etc. În acest mod, problema se reduce la calcularea frecvenței pe parcursul unui interval constant de timp, având întotdeauna ca final ziua curentă. Cineva ar putea realiza acest lucru prin simpla căutare a cuvântului în depozit și numărarea aparițiilor care se încadrează în intervalul de referință – o funcție care poate fi apelată doar o dată într-un lung interval –

să zicem de la doi la cinci ani (deoarece nu ne putem aștepta ca marcarea „învechit” să se apară brusc, de ieri până astăzi, de exemplu ...).

Este clar ca orice decizie finală cu privire la aceste poziții trebuie să fie luată de o autoritate științifică (în speță, Academia). Deciziile acesteia însă trebuie să se bazeze pe semnalele trimise de *Portal*, care la rândul lor trebuie să aibă la bază probe statistice clare.

Alte fluxuri de procesare pot implementa alte funcționalități. Un număr de resurse în format electronic, care sunt de o importanță majoră în menținerea unei limbi actualizate din punct de vedere tehnologic, pot fi conectate permanent la *Portal*. Unele dintre acestea ar putea fi: Dicționarul Tezaur oficial al limbii (eDTLR în cazul limbii române, vezi (Cristea et al., 2009)), WordNet (Fellbaum, 1998), VerbNet (Kipper et al., 2008), FrameNet (Fillmore, 1976; Atkins et al., 2003) – pentru a numi doar câteva. Presupunând ca aceste resurse devin, la un anumit moment, complete pentru limba *L*, ele trebuie actualizate permanent cu evoluția limbii. Astfel, dinamica limbii trebuie să se reflecte și în resurse. Dacă, așa cum se sugera mai sus, fiecare element lexical are o fișă proprie pe *Portal*, atunci ea ar trebui să includă referințe în toate resursele. Ca atare, fișa cuvântului *w* este legată la intrarea lui din Dicționar, care conține și un inventar al sensurilor lui, iar aceste sensuri sunt sincronizate cu cele conținute în WordNet pentru acest element lexical, precum și cu intrarea sa din VerbNet și FrameNet. Toate aceste resurse sunt conectate între ele și sunt actualizate în ritm cu evoluția limbii de către *Portal*.

Portalul poate, de asemenea, găzdui un număr de servicii adresate *furnizorilor de resurse*, cercetătorilor limbajului, consumatorilor industriali ori publicului larg. Serviciile publice ori comerciale pot fi supuse unor taxe, profitul urmând a fi întors *furnizorilor de resurse*, de exemplu proporțional cu contribuția lunară a fiecăruia pe *Portal* (măsurată în număr de caractere) ori cu numărul de accese ale consumatorilor la propriile resurse.

Pot fi imaginate și alte tipuri de servicii plătite, cu direcționarea profitului către *furnizorii de resurse*. Un exemplu îl constituie postări care fac reclamă cărților tipărite, acces online la fragmente ale publicațiilor ce se găsesc pe piață etc. Cred că posibilitățile de servicii ce ar putea fi oferite de *Portal* prin aplicarea unor rafinate tehnologii lingvistice și pe baza cărora *furnizorii de resurse* să obțină profit sunt extrem de diverse și ar trebui promovate cu insistență. Acest lucru ar putea diminua teama *furnizorilor de resurse* într-o lege care, doar aparent, afectează controlul lor asupra proprietății, prin aceasta, contribuind la micșorarea rezistenței lor față de aplicarea ei în practică.

5. Evaluare

Este clar că o astfel de inițiativă va aduce zilnic pe *Portal* un volum foarte mare de date lingvistice. O estimare, chiar și grosieră, a necesităților de procesare și a costurilor presupuse de un astfel de demers la nivel național ar trebui să ia în considerare parametri cum ar fi: numărul de edituri înregistrate, numărul mediu de publicații per editură per an, numărul mediu de pagini a unui obiect tipărit, numărul mediu de caractere pe pagină. Lăsând la o parte publicațiile episodice de mici dimensiuni, ancheta

noastră cu privire la cantitatea medie de date publicate în cărți și periodice, într-o țară de mărime medie din Europa cum e România, la nivelul anului 2008, a evidențiat un volum de date în format text care nu atinge 1 GB pe zi.

Un canal cu lățimea de bandă de 12.5 Mb/sec poate asigura cu ușurință transferul necesar descris în secțiunea 3, evitând blocajele chiar și în momentele de supraaglomerare. Pentru motive de siguranță ar trebuie asigurată o încărcare distribuită iar datele ar trebui stocate redundant (*mirroring*) în cel puțin două centre, în locații diferite. Cum au dovedit deja giganții Internetului (de exemplu: Google⁸), tehnologia RAID, care grupează un număr foarte mare de mici calculatoare, este o soluție ieftină și care ar putea fi avută în vedere pentru a asigura conservarea pe termen lung și pentru obținerea unei viteze de procesare confortabilă.

6. Concluzii

Avantajele unei infrastructuri dedicate prelucrării datelor lingvistice, de o anvergură atât de mare ca cea schițată în această lucrare, sprijinită și de o legislație favorabilă, sunt greu de stabilit corect acum. În primul rând, ea oferă o soluție completă pe termen lung pentru conservarea datelor lingvistice ale unei națiuni, precum și o radiografie aproape completă a evoluției ei diacronice. În al doilea rând, ea pune bazele unor cercetări exhaustive asupra limbii. În al treilea rând, ea aduce în atenție o gamă largă de aplicații atrăgătoare din punct de vedere comercial, spre beneficiul autorilor de texte, *furnizorilor de resurse*, și consumatorilor de limbă, prin aceasta accelerând cercetările din domeniul lingvisticii computaționale și al prelucrării limbajului natural și influențând benefic nivelul de cultură al maselor.

Succesul unei astfel de inițiative la nivel național depinde foarte mult de o viziune europeană concertată. Suflul nou care se simte la momentul actual în Europa cu privire la realizarea de infrastructuri pentru procesare lingvistică, la stabilirea de standarde pentru reprezentarea datelor lingvistice, la achiziționarea de resurse lingvistice și susținerea financiară comunitară a unor proiecte transnaționale precum CLARIN⁹, FlareNet¹⁰, T4Me¹¹, Meta-Net¹² etc. ar trebui să creeze și premisele unor legislații favorabile. Propunerea avansată în această lucrare este conformă cu alte inițiative care încearcă să crească nivelul de conștientizare cu privire la necesitatea accesului liber la știință. În nici un caz lucrarea nu pledează împotriva proprietății intelectuale (Stephan, 2001), dar este în favoarea reconsiderării legislației cu privire la drepturile de autor, care, în forma actuală, este extrem de restrictivă în multe cazuri relativ la utilizarea resurselor lingvistice în folosul cercetării. La urma urmei, limba noastră, așa cum o folosim astăzi, este creația colectivă a națiunii și este supusă unor transformări (n-aș zice că acestea ar putea fi caracterizate întotdeauna drept cizelări...) permanente, realizate de către toți vorbitorii ei. Donarea creațiilor lingvistice pentru conservarea limbii și cercetare, nefiind în nici un fel dăunătoare pentru creator, nici intelectual și nici material, reprezintă, în definitiv, minimum de recompensă pe care un autor care

⁸ <http://infolab.stanford.edu/~backrub/google.html>

⁹ www.clarin.eu

¹⁰ <http://www.flarenet.eu/>

¹¹ <http://t4me.dfki.de/>

¹² <http://www.meta-net.eu/>

utilizează limba îl datorează celor care au inventat-o, în beneficiul celor care o vor folosi în viitor.

Referințe bibliografice

- Atkins, S., Rundell, M. and Sato, H. (2003). The Contribution of Framenet to Practical Lexicography, *International Journal of Lexicography*, Volume 16.3: 333-357.
- COM (2009) 532 – Communication from the Commission. Copyright in the Knowledge Economy.
- Cristea, D., Răschip, M., Moruz, A. (2009). Steps in Building the Electronic Version of the Thesaurus Dictionary of the Romanian Language, in Proceedings of the IVth National Conference The Academic Days of the Academy of Technical Science of Romania, ASTR – Filiala Iași și Universitatea Tehnică „Gheorghe Asachi” Iași, Ed. Agir, ISSN 2006-6586.
- Fellbaum, C. (1998). WordNet: An Electronic Lexical Database, MIT Press.
- Fillmore, C. J. (1976). Frame semantics and the nature of language. In *Annals of the New York Academy of Sciences: Conference on the Origin and Development of Language and Speech*, Volume 280: 20-32.
- Kipper, K., Korhonen, A., Ryant, N., Palmer, M. (2008). A Large-scale Classification of English Verbs, *Language Resources and Evaluation Journal*, 42(1), pp. 21-40, Springer Netherland.
- Kunze, J. and R.P.C.Rogers (2003). The ARK Persistent Identifier Scheme. Internet draft at <http://www.cdlib.org/inside/diglib/ark/arkspec.pdf>.
- Schwardmann, U. (2009). PID System for eResearch. EPIC – the European Persistent Identifier Consortium, personal communication at NEERI-09, Helsinki.
- Stephan, K. (2001). Against Intellectual Property. *Journal of Libertarian Studies*, 15.2, 1-53.

COMUNICAREA ELECTRONICĂ ȘI PROBLEMELE NOASTRE ORTOGRAFICE – FĂRĂ SOLUȚII?

LUCIAN CHIȘU

Muzeul Național al Literaturii Române, București

lucianchisu@gmail.com

Rezumat

Acest text poate fi considerat o continuare a unei intervenții anterioare a autorului, apărute sub titlul *Comunicarea electronică și problemele noastre ortografice* (AUSH, seria *Jurnalism* 2003). La acel moment, limbajul informatizat în limba română prezenta un aspect mai degrabă tranzitoriu. Autorul considera că, în scopul bunei sale utilizări, ar fi fost necesare două tipuri de intervenții. Pe de o parte cea a specialiștilor IT în scopul realizării softului și adaptării claviaturii la limba română, iar pe de alta, atitudinea fermă a specialiștilor (lingviști, filologi, elite culturale) la respectarea normelor ortografice și de limbă literară și în cadrul limbajului informatizat, pe cale de a se generaliza în ceea ce privește comunicarea. Concluzia degajată din noul conținut, incluzând perioada 2003-2010 este următoarea: deși există soft-uri și sunt comercializate claviaturi românești, deși s-au promulgat acte normative care prevăd sancțiuni în cazul nerespectării limbii literare, abaterile, greșelile și ezitățile de la începuturi se manifestă într-o proporție îngrijorătoare. Sunt indicate site-uri ale oficialităților, unele aparținând Ministerului Culturii și Ministerului Educației, în care, departe de a fi puse în practică, reglementările ortografice sunt ocolite, în unele cazuri cu indicația expresă de a se evita diacriticele.

I. Începând cu ultimul deceniu al secolului trecut, utilizarea în domeniul presei a limbajului informatic și-a făcut simțite avantajele, răsfrânse la scurt timp și în cultura română. Pe de o parte, creșterea fără precedent a gamei și numărului de servicii electronice (sub formă textuală, audio, grafică și/sau video), iar pe de alta, accesul tot mai nelimitat, inclusiv cel individual și mai ales al tinerilor, au creat premisele dezvoltării exponențiale ale noilor forme adaptate comunicării. În primii ani ai mileniului al treilea fenomenul se generalizase, devansând tipurile de comunicare lingvistică până nu demult consacrate. Totuși, chiar și cei mai fideli dintre susținătorii informatizării limbajului erau de acord cu privire la imensa cantitate de necunoscut care însoțea procesele fenomenului în discuție¹.

¹ „În această epocă năvalnică de schimbări tehnologice rapide, toți luptăm pentru a ne păstra propria direcție. Progresele care se dezvoltă zilnic în comunicații și tehnică de calcul pot fi înspăimântătoare și deusolante. Este firesc, în contextul dat, să ne întrebăm dacă aceste schimbări sunt bune sau rele, dacă trebuie întâmpinate cu bucurie sau cu teamă. Răspunsul este: și una și alta.” (*Manifestul tehnorealist*, 2000, p. 307). Textul reprezintă rodul colaborării a 12 scriitori preocupați de tehnologie: David Bennahum (editor al publicației „Meme”), Brooke Shelby Biggs (editorialist la „San Francisco Bay Guardian”, varianta electronică), Paulina Borsook (scriitor), Marisa Bowe (fost redactor-șef al revistei „Word”), Simson Grafinkel (colaborator la „Wired”, editorialist la „The Boston Globe”), Steven Johnson (scriitor, redactor-șef la „FEED”), Douglas Rushkoff (scriitor, editorialist la „Time Digital”), Andrew Shapiro (fellow la Centrul Berkman pentru Internet și Societate al Facultății de Drept Harvard, colaborator la „The Nation”), David Shenk (scriitor, comentator la The National Public Radio), Steve Silberman (editor în probleme de cultură la „Wired”), Mark Shtalman (scriitor, cofondator al „FEED”). Conceptul acestui document a fost schițat de

1. Deoarece în anii imediat următori momentului utilizării computerului, în țara noastră nu existau „tastaturi” cu claviatura adaptată la alfabetul limbii române, apăruseră adevărate modele de scriere constând în abateri de la normele ortografice. Alături de texte redactate fără semnele grafice (diacritice) ale literelor *ă, â, î, ț, ș*, al căror înțeles deplin rezulta contextual, în tot mai extinsa comunicare pe suport electronic se dezvoltau noi tipuri de scriere în care alfabetul românesc era folosit *ad libitum*, după priceperea sau inventivitatea utilizatorilor. Chestiunea rigorii în ceea ce privește recuperarea sensului cuvintelor devenise pentru unii dintre utilizatori stringentă. Sensul lexical trebuia desprins din context pentru că, din cauza absenței diacriticelor, nu rezulta clar dacă în enunțul respectiv se făcea referire la *bulgări* (de pământ ori de zăpadă) sau *bulgari*, cetățeni ai țării vecine („*loviti de bulgari*”), de *paturi* sau *pături* (*sinistratilor li s-au oferit paturi*), dacă în *supa de cocos*, produsul culinar provenea din carnea de *cocoș* sau din fructul exotic *cocos*, care intră ca ingredient alimentar în compoziții asemănătoare. În pofida acestor neajunsuri, abaterile de la norme își croiseră calea și sub alte forme, la fel de insolite². Metodele de aplicație fiind iluzorii, ele înseși constituiau erori sau devieri de la normă și ilustrau totodată accesul dificil, greoi, complicat și inutil la un sistem (tastatură) impropriu.

Lucrul se datora în primul rând unei carențe legislative, dar, în ansamblu, fenomenul prezenta o dublă cauză.

a. Cea dintâi, de natură tehnică, necesita realizarea softului și a tastaturilor pentru limba română, fiindcă inițial cele comercializate în România fuseseră proiectate pentru uzul în limba engleză. Sarcina realizării „programelor” pentru diacritice și a tastelor aferente alfabetului românesc cădeau exclusiv în seama programatorilor IT și inginerilor de sistem din țara noastră. Chestiunea a fost întrutotul depășită și a devenit de domeniul trecutului.

b. Consecință nemijlocită a celei dintâi, a doua cauză servea pragmatismului greșit înțeles: utilizarea corectă a limbii române în scrierea electronică trecuse pe un al doilea plan, cedând întâietate beneficiilor, fapt care a favorizat instalarea aproape nestânjenită a unui autentic haos ortografic. Contau mai mult imensele avantaje ale scrierii electronice constând în fiabilitate, actualizare, conservare, multiplicare, dar mai ales eficiența, rapiditatea și confortul comunicării electronice, spațiul și timpul fiind contrase în ceea ce numim oarecum impropriu „mediul virtual”. Cu alte cuvinte, comunicarea electronică se afla în coliziune cu alfabetul și ortografia noastră normativă, situându-se, spre nemulțumirea unora dintre utilizatori, într-un punct tranzitoriu de care nimeni nu părea a fi responsabil să ia atitudine. S-a instalat, nu fără temei, teoria formelor fără fond, valabilă în tipuri de situații în care „forțele reacționare”, înșelate în așteptările lor, declanșează conflictul dintre procesul vieții (care nu știe ce-i răgazul) și formele de exprimare culturală (care nu știu ce-i schimbarea). Sentimentul nemulțumirii se conturează în astfel de situații sub forma crizei.

Shapiro, Shenk și Johnson. A fost publicat pe 12 martie 1998. Poate fi citit în versiune originală la www.techmorealism.org. Vezi, Ana Tăbîrcă, traducătoarea *Manifestului tehmorealist* (***, 2000).

² Exemple privind improprietățile semantice ale textelor scrise fără diacritice, cu mult mai numeroase, care proliferaseră în limbajul informatic al anilor 2000, am oferit în (Chișu, 2002) și în studiul din (Chișu, 2003).

După trecerea unui număr de ani, dezavantajele aspectului prim (tehnic) au fost anulate, însă consecințele începuturilor n-au dispărut, putându-se afirma că „problema”, deși părea a fi fost rezolvată, încă persistă.

În acest context ambiguu, singura certitudine o reprezintă generalizarea scrierii pe suport electronic, care se manifestă ca un fenomen (inter)național tocmai ca urmare a imenselor foloase pe care le generează. Ca pretutindeni pe glob, instituțiile oficiale și private din România, precum și uriașa masă de utilizatori individuali au trecut la noul format. Indiferent de caracterul oficial sau particular al textelor, majoritatea covârșitoare a lor este procesată electronic și abia după aceea transferată pe suportul tradițional (hârtia). Utilizăm fără îndoială, la scară națională, limbajul informatic cu extraordinarele sale disponibilități și, în continuare, însoțite de „reflexele” vechilor abilități.

2. Este, iarăși, o certitudine că, în privința limbilor „culte”, limbajul informatic reprezintă tehnic cea mai spectaculoasă ascensiune. Dar, în ceea ce privește culturile, ale căror sisteme lingvistice includ în evoluția lor alte semne grafice decât cele consacrate pentru limba engleză, limbajul informatic a exercitat, cel puțin pentru o perioadă de timp, un altfel de impact. Limbajul constituie fundamentul comunicării între oameni și în mod indiscutabil pentru toți aceștia poartă profunde conotații emoționale și culturale provenind din moștenirea lor literară, istorică, filosofică și educațională. De aceea, limba maternă nu trebuie să reprezinte un obstacol în calea accesului la cunoașterea multiculturală umană disponibilă în noul mediu. Dacă utilizarea ei ortografică devine, cel puțin pentru unele semne grafice, aleatorii, obstacolele nu întârzie să apară.

II. Cu toate că, pe fond, chestiunea utilizării limbajului informatic în spațiul cultural românesc și-a găsit soluția optimă, între altele și ca urmare a emiterii și promulgării unor legi (Hotărâri de Guvern), numeroși dintre utilizatorii săi împărtășesc sentimentul că de la soluție la rezolvarea celui de-al doilea aspect, cel practic, rezistența în fața regulilor, normelor și corectitudinii continuă să se manifeste în proporții de masă, chiar dacă efectele recuperatorii își vor face simțite în cele din urmă, cândva, consecințele. Aceștia nu caută vinovați sau vinovații dar, pe de altă parte, nici nu pot accepta că s-a asistat (pasiv) la un fenomen... natural, deci greu de stăpânit, incontrollabil. În toți acești ani de evoluție nestânjenită a unui fenomen care devenea pe zi ce trece previzibil și ireversibil o nouă formă generalizată de comunicare, ar fi fost de așteptat reacții și atitudini din rândul specialiștilor limbii. Punctul de vedere al acestora asupra respectării normelor ortografice și, implicit al limbii noastre literare în comunicarea electronică ar fi avut cea mai mare greutate și autoritate. În acest interval, reacția dumnealor a fost foarte puțin sau deloc vizibilă, cei mai mulți resemnându-se a consemna adevărul că limba română se află într-o grea suferință. Totuși, domniilor lor, doctorilor în... litere le-ar fi revenit un rol mai activ.

1. Era, de asemenea, datoria tuturor intelectualilor, de toate formațiile umaniste ori scientiste, să ia act de autentică mutație intervenită, care favorizează saltul în SI-SC și în viitorul comun. Această formă de sprijin ar fi constituit cel mai important element de solidarizare cu specialiștii în limbaje (lingviști, filologi) care, altfel, împart o soartă ingrată, știut fiind că primii dintre aceștia studiază mai întâi evoluțiile și tendințele limbajului, mereu schimbător, spre a formula eventual recomandări și finalmente

norme, respectate cu sfințenie de filologi și nu numai de ei, ci de toți indivizii care formează societatea.

a. Trebuie ținut cont, de asemenea, că promovarea limbii române în SI-SC sub aspectele sale scris-vorbit a devenit, cu un cuvânt ale cărui semnificații sunt mai degrabă arhaice, apanajul altor domenii de investigare, în sensul acesta vorbindu-se curent despre „ingineria limbajului”, aplicație care vine să sublinieze aspectele legate de validarea experimentală a ipotezelor științifice prin modele riguros verificate. În format electronic, limbajul scris – vorbit ocupă numai o avanscenă (de resursă lingvistică), domeniul de cercetare devenind prelucrarea automată a limbajului natural. Acesta din urmă reflectă progresele științifice și tehnologice care-l desprind de momentul inițial (lingvistica formală) spre a accede în zonele generării de modele abstracte (lingvistica matematică) și spre lingvistica computațională, constând în prelucrarea limbajului natural ca *summum* al teoriilor și modelelor de inteligență artificială. Aceste forme de cunoaștere a limbajului abordează problematica de comunicare om-calculator.

b. Este, de asemenea, demn de consemnat faptul că, încă din primii ani de manifestare a limbajului informatic, specialiștii noilor arii de cercetare, au ținut să accentueze ideea că dezvoltarea armonioasă a societății informaționale bazată pe cunoștințe este posibilă doar prin promovarea informației și accesului cu caracter multilingv și multicultural. Mai mult chiar, în raportul special al comisiei Europene (*Towards a European Language Infrastructure*), raportorul A. Danzin atrăgea atenția asupra următorului pericol: „În era electronică, este esențial pentru supraviețuirea unei limbi ca ea să fie folosită în sistemul de informare electronică”, rezultând cel puțin indirect că procesul de informatizare a unei limbi naturale nu înseamnă în niciun chip stâlcirea ei”.

Sub aceste aspecte, datoria celor mai responsabili dintre utilizatorii limbii noastre rămăsese doar un deziderat.

III. La cele afirmate până acum este necesar să adăugăm o serie de aspecte particulare ce se regăsesc în evoluția istorică a comunicării scrise în limba română.

1. ÎNSEȘI „realitățile” interne, istorice, ale scrierii în limba română, au cunoscut etape punctate de continue dispute și neînțelegeri de-a lungul epocilor. În ceea ce privește „biruința scrisului în limba română”, înaintașii noștri s-au servit inițial de un alfabet cu grafie slavonă de 43 de slove (alfabetul chirilic). Concomitent cu utilizarea greoiului alfabet chirilic, întrupând „sufletul românesc în veșmânt slavon” (cum l-a numit Nicolae Cartoian), la sfârșitul secolului al XVIII-lea, învățații ardeleni se foloseau de grafia latină³, dar intenția reprezentanților *Școlii ardeleni* se îndrepta spre utilizarea unui alfabet etimologic, care ar fi dus la o supărătoare ruptură între modul de a scrie și pronunța cuvintele. Astfel de fapte erau și mai mult adâncite de caracterul incipient, pentru acele vremuri, al utilizării grafiei. „Regulile ortografice propuse în urma discuțiilor derulate între 1866 și 1881 de Academia Română întâmpinau dificultăți deosebit de mari în răspândirea lor și rezistență pe măsura dificultății lor” (Zugun, 1992).

³ Se știe, spre exemplu, că Ion Budai-Deleanu a optat primul pentru grafiile ț și ș, așa cum sunt redată astăzi.

Considerând că au fost schițate câteva elemente relevante pentru începuturile alfabetului și evoluției ortografiei la noi, informații care se regăsesc în studii și cărți publicate de numeroși specialiști (Avram, 1990), (Beldescu, 1985), (Hristea, 1981), (Nica și Cureteanu, 1980), (Jordan, 1978), (Șuteu, 1976), (Graur, 1974), (Ciobanu și Sfârlea 1970), (Rosetti, 1967), consemnăm că discuții polemice, urmate de propuneri, reforme și hotărâri ortografice au avut loc în anii 1881, 1904, 1925, 1932, 1953, 1957 și 1991, acestea nestingându-se⁴ niciodată.

Reținem însă, că în atmosfera tensionată, dacă nu uneori incendiară, a acestor dispute, totuși s-a convenit asupra necesității unui mai pronunțat caracter coercitiv al normelor ortografice, după cum rezultă și din unele nu foarte îndepărtate în timp luări de poziție în istoricul problemei.

2. La avatarurile alfabetului și normelor ortografice succint enunțate mai înainte, se adaugă hotărârea academică ratificată în ședința Adunării Generale a Academiei Române din 31 ianuarie 1991, a cărei contestare a avut drept efect scindarea alfabetului. Numeroase edituri și publicații, unele de prestigiu, au ignorat hotărârea în discuție⁵.

În plus, concomitent cu „schisma alfabetică” din 1991, odată cu pătrunderea masivă în societatea românească a tehnologiilor multimedia aveau să se ivească noi situații conflictuale. În uz s-a instalat și cea de a treia formă, scrierea electronică fără diacritice, cu variantele ei opționale. Ideea eludării convenției și ignorarea efectului normativ a generat frecvent puncte de vedere orgolioase, unul dintre acestea fiind exprimat de un cunoscut critic literar, universitar și academician, director al celei mai cunoscute publicații literare românești: „Iar dacă dorim cu tot dinadinsul să facem convenția cât mai firesc-utilă, atunci, date fiind schimbările produse în scriere de calculatoarele moderne [noi am fi spus *de utilizatorii calculatoarelor*], ar fi poate necesar să abandonăm sedilele de la *ț* și *ș* și să scriem *tz* și *sh*, [internaționalizând] cumva scrierea, sau chiar să scriem, ca portughezii, pe *ă* ca *a*, indiferent de pronunție, și pe *î* ca *i* (cum de altfel calculatorul o și face!) [*noi am sublinia că nu calculatorul, ci utilizatorii o ...fac!*]. Dar oare reformiștii noștri vor fi de acord? Nu cred, fiindcă problema lor nu e de a avea o convenție care să fie lesne respectată, ci, una, cum să zic? *respectabilă*”

⁴ În 1881, prin contribuția decisivă a junimiștilor, reveniți la Academie după slăbirea inevitabilă a poziției, dominante până atunci, a latiniștilor, a fost salutară impunerea ortoepiei limbii, prin influența scrierii fonetice asupra vorbirii, putându-se crea, numai astfel, o benefică circularitate sistemică *vorbită îngrijită — scriere fonetică — ortoepie*, imposibilă în cazul unei ortografii bazate pe etimologism. Între 1881- 1904 s-a realizat trecerea de la etimologism la fonetism, etapă de tranziție ce „a întărit unele modalități de scriere care s-au impus, ca, de exemplu, scrierea vocalei [ă] prin a cu un semn diacritic deasupra, a consoanelor [ș] și [ț], prin ș și respectiv ț, cu sedile ș.a.m.d.” (apud Petru Zugun, 1990, p. 42). Un alt proiect academic a fost redactat în 1932 de Ov. Densusianu. El a avut drept moment inițial discuțiile începute încă din 1925 de Sextil Pușcariu, reprezentând modificări sau precizări aduse reformei ortografice din 1904. O nouă simplificare și clarificare a ortografiei românești realizează Academia în 1953 cu importanta completare din 1965, când este normată, tot de Academie, folosirea literei *â* în familia cuvântului *român*. În anul 1991 Academia Română, prin președintele ei de atunci, dl. acad. Mihai Drăgănescu, a luat hotărârea de a recomanda câteva amendamente la normele ortografice ale limbii române din 1953: a) utilizarea lui *â* în locul literei *î* în interiorul cuvintelor; păstrarea literei *î* în interiorul cuvintelor dacă acestea sunt rezultatul unui cuvânt compus, în care cel de-al doilea element lexical îl conține pe *î* inițial; b) utilizarea formei *sunt* în locul formei *sînt*, procedându-se în mod similar cu formele flexionare. Această hotărâre ratificată în ședința Adunării Generale a Academiei Române din 31 ianuarie 1991, a stârnit vii proteste și noi puncte de vedere polemice. În aceeași perioadă apare și se dezvoltă exponențial scrierea electronică, mijlocită de computer.

⁵ Astăzi avem situația, cel puțin bizară, a editurilor care, atunci când tipăresc manuale, școlare folosesc ortografia academică, iar când editează alte cărți o utilizează pe aceea ...contestată.

(Manolescu, 2002). Cei de la care era de așteptat „tăierea” acestui veritabil nod gordian nu s-au pronunțat, dar acest derapaj n-a fost deloc pe placul informaticienilor.⁶

3. Iată cum se (re)prezintă limba română în format electronic pe unele dintre site-urile oficiale, la ora de față:

a. CNCISIS (Consiliul Național al Cercetării Științifice din Învățământul Superior), organ consultativ al Ministrului Educației, Cercetării, Tineretului și Sportului, exprimând totodată punctul de vedere al comunității științifice în ceea ce privește politica cercetării științifice apare postat fără diacritice. Unele fișiere sunt dedicate consultării în limba română, altele în limba engleză, dar între acestea nu există practic nicio diferență. La rubrica *Instrucțiunile privind completarea formularului „cerere de finanțare”* se face precizarea: „este OBLIGATORIU ca în timpul completării formularului să nu introduceți caracterul ghilimele, caractere diacritice”. Tot astfel stau lucrurile, adică fără prezența diacriticelor, și cu CNFIS (Consiliul Național pentru Finanțarea Învățământului Superior), UEFISCSU (Unitatea Executivă pentru finanțarea Învățământului Superior și a Cercetării Științifice Universitare), parțial cu AFCN (Administrația Fondului Cultural Național), care se subordonează Ministerului Culturii și Patrimoniului Național, unde, în cadrul rubricii „întrebări frecvente”, se recomandă evitarea folosirii normelor obligatorii de scriere, cu diacritice.

b. Notariatele folosesc în documentele legalizate diacriticele numai la cererea expresă a clientului. În celelalte situații, formula fără diacritice este generalizată, în pofida faptului că pot apărea confuzii supărătoare, mai ales de nume.

c. În marea majoritate a Universităților din țara noastră sunt acceptate Teze de licență și de masterat redactate în mod frecvent fără diacritice. S-au întâlnit cazuri, însoțite de naive surprinderi, că ar fi greșită redactarea fără diacritice, inclusiv la limba și literatura română. Faptul nu este în măsură să enerveze, ci, dimpotrivă, să se generalizeze.

d. Cele mai multe dintre site-urile de pe net perpetuează fără jenă situația descrisă mai sus.

Exemplele ar putea continua dar, din cele prezentate, se detașează următoarea observație: raportul dintre scrierea în limba română literară și celelalte formate înclină majoritar în favoarea celui nerecomandabil, greșit. Singura concluzie, din păcate „viabilă”, se referă la adevărul că generalizarea scrierii electronice a devenit o realitate de necontestat.

4. Dacă nu vom trece sub tăcere că acțiunea globalizării are consecințe devastatoare asupra culturilor identitare, atunci trebuie să invocăm și studiul lui (Bârlea, 2009) care demonstrează indubitabil, prin intermediul unor cercetări efectuate de specialiști ai Comunității Europene, că dispariția repercusiunile asupra limbilor

⁶ „Există o sumedenie de situri unde limba română apare trunchiat, adică extirpată de diacritice. Prima explicație care ți se servește este problema anumitor greutăți dificil de contorizat. Dacă acum ceva timp exista un sâmbure de adevăr, acum e foarte greu de susținut o astfel de explicație, nocivă după părerea mea. Cu mult mai grav e că există profesori în învățământul românesc care în loc să promoveze corecta folosire bat apa-n piuă cu subterfugii. Dar nu numai aici vezi astfel de «scăpări». Există o serie de instituții culturale, edituri, reviste de exemplu, unde o astfel de explicație își face loc cu coatele. Am fost surprins să aud din gura unui critic literar și conducător de revistă o astfel de ineptie. La contra argumentele mele, replica a fost un soi de făcut cu mâna din categoria «merge și așa» (Dan Iancu, redactor șef „PCMagazine România”, p. 4.).

naționale reprezintă un fenomen real și destul de amenințător, ca și topirea calotei glaciare. Printre cauzele dispariției limbilor, autorul articolului citat enumeră: reacția tinerelor generații față de limba maternă, potențialul social-economic al comunităților de vorbitori, relația scris-oral, politicile oficiale. Desigur, faptele prezentate mai înainte nu reprezintă, în acest moment, un pericol real pentru limba română. Din studiul cercetătorului înainte amintit, rezultă, pe bază de date și informații, că la fiecare două săptămâni dispăre câte o limbă având mai puțin de 200.000 de vorbitori. În acest ritm, accelerat de noile media, ștergerea identităților se va produce în următorii o sută de ani. E bine? E rău? Apariția unui *model unic*, globalizant, pentru noile situații create, generează întrebări și solicită răspunsuri la interogații diverse. Nu se schimbă, prin urmare, funcțiile sau domeniul comunicării prin limbaj, ci, prin noile tehnologii instalate, pe lângă imense avantaje, se profilează și un nou mediu și tip de cultură, aflat în conflict nedecarat cu vechile modalități de comunicare tradițională, responsabile cu metabolismul organismelor culturale. E bine? E rău?

Ca instrument (proteză) a informației și cunoașterii, în pofida faptului că a pătruns adânc în evoluția societății, în țara noastră nu a existat și nu există o politică de stat coerentă cu privire la această mare invenție, cea mai discutabilă dintre problemele ridicate în discuție fiind aceea a patrimoniului național care deține un specific național sau identitar puternic conturat.

IV. Concludem, așadar, că, privitor la dezbaterile asupra scrierii, inclusiv ale celei pe suport electronic, în lumea filologilor formula *grammatici certant*, își pune amprenta, etern actuală. Prin urmare, în absența unor decizii ferme, deși instituționalizate, s-a ajuns la situații realmente de neînțeles.

Programul conferinței *ConsILR 2010 (Resurse lingvistice și instrumente de lucru pentru prelucrarea limbii române)*, în cadrul căruia își unesc forțele informaticieni, cercetătorii în inteligența artificială și lingviști se situează la polul diametral opus. Chestiunile supuse dezbaterii privesc *realizarea de resurse lingvistice românești, textuale ori vorbite, în forma originală ori adnotată, crearea de corpusuri românești reprezentative, realizarea de colecții lexicografice românești în format electronic, tehnologii lingvistice aplicate limbii române, aplicații în care au fost utilizate tehnologii lingvistice pentru limba română, realizări de lingvistică teoretică cu aplicații în tehnologia limbii române, proiecte de cercetare ce implică dezvoltarea de resurse și instrumente dedicate limbii române.*

Față de cele prezentate în această intervenție, *Conferința* găzduită de Muzeul Național al Literaturii Române devine un spațiu aulic, locul perfect al schimbului de idei pentru viitorul Si-SC. Și, totuși, rămâne sentimentul extrem de disconfortant că între realitatea palpabilă a utilizării limbii române în comunicarea de tip electronic și obiectivele *Conferinței*, repet de o utilitate maximă, se interpune, metaforic vorbind, povestea episodului căderii orașului Constantinopol, când, în timp ce fortăreața era escaladată de păgâni, în cercurile înaltelor spirite (religioase) din acel loc se purta o discuție extrem de animată despre ...sexul îngerilor. Desigur, afirmația trebuie luată numai ca o metaforă, poate nereușită, dar căreia am încercat să-i aducem câteva argumente, după opinia noastră, indubitabile.

Tocmai de aceea, se poate admite că, paralel cu aceste inițiative atât de demne de interes, aspectele „practice” ale uzului limbajului în format electronic, semnalate pe parcursul acestei intervenții, merită readuse în actualitate, în fața participanților la *ConsILR 2010*.

Referințe bibliografice

- *** (2000), Manifestul tehnorealist, în „Secolul XX”, nr. 4-9, 2000 (421-426), pp. 307-311 (traducerea Ana Tăbîrcă).
- Avram, M. (1990), Ortografie pentru toți (30 de dificultăți), Editura Academiei Române, București.
- Bârlea, P. G. (2009), Dispariția limbilor - catastrofă umană sau formă naturală de schimbare lingvistică?, în „Caiete critice”, nr. 8-9 (262-263)/2009, pp. 97-117.
- Beldescu, G. (1985), Ortografia actuală a limbii române, Editura Științifică și Enciclopedică, București.
- Chișu, L. (2002), O chestiune urgentă, în „Tribuna învățământului”, an LIII, nr. 634 (2515), 18-24 martie 2002, p. 1; 2.
- Chișu, L. (2003), Comunicarea electronică și problemele ortografiei noastre, în „Analele USH”, Seria Jurnalism, IV, nr.4, pp. 31-48.
- Ciobanu, F., Sfirlea, L. (1970), Cum scriem, cum pronunțăm corect (norme și exerciții), Editura Științifică, București.
- Graur, Al. (1974), Mic tratat de ortografie, Editura Științifică, București.
- Hristea, T. (1981), Sinteze de limba română, ediția a II-a, revăzută și mult adăugită, Editura Didactică și Pedagogică, București.
- Iordan, I. (1978), Istoria lingvisticii românești, Editura Științifică și Enciclopedică, București.
- Manolescu, N. (2002), Cum scriem, în „România literară, An XXXV, nr. 38, 18-24 septembrie 2002, p. 1.
- Nica, M., Cureteanu, S. (1980), Predarea ortografiei în gimnaziu, Editura Didactică și Pedagogică, București.
- Rosetti, Al. (1967), Introducere în fonetică, Editura Științifică, București.
- Șuteu, F. (1976), Influența ortografiei asupra pronunțării literare românești, Editura Academiei, București.
- Zugun, P. (1992), Ortografia limbii române. Trecut, prezent, viitor, Institutul European, Iași, 1992, p. 40.

FACTORII CARE INFLUENȚEAZĂ MOMENTUL OPTIM DE MIGRARE LA DIACRITICELE CORECTE ÎN LIMBA ROMÂNĂ¹

BOGDAN STĂNCESCU

S.C. Moongate Video Production srl, București – România

bogdan@moongate.ro

Rezumat

Preconizez că în perioada 2010-2015 va avea loc în cea mai mare parte migrarea conținutului de limbă română de la utilizarea caracterelor cu diacritice cu sedilă (conform ISO-8859-2) la caracterele cu diacritice corecte cu virgulă (conform Unicode 3.0). Acest document încearcă o sinteză a factorilor care influențează, de la caz la caz, momentul ideal de migrare.

1. Introducere

Semnele diacritice din partea de jos a caracterelor românești „ș” și „ț” sunt virgule. Acest detaliu este de la sine înțeles pentru orice vorbitor nativ de limbă română: este un fapt nedisputat care se învață în clasele primare și nu mai trebuie repetat niciodată în mod explicit. Iar asta într-o asemenea măsură încât însăși Academia Română nu a simțit nevoia să se pronunțe în această privință decât în anul 2003, și chiar și atunci numai pentru că a răspuns unei întrebări explicite în acest sens.²

Atunci când a fost creată prima codare a caracterelor pentru Europa de Est, în 1987 (Latin-2³), caracterele pentru limba română au fost comasate cu cele pentru alte limbi din această zonă geografică scrise în mod uzual cu grafie latină, precum ceha, maghiara, poloneza și altele. Între limbile asociate acestui standard, limba română este singura care folosește caracterele „ș” și „ț”, sau orice caractere similare din punct de vedere vizual.

Caracterul „ș” din limba română („s cu virgulă”) este foarte similar din punct de vedere vizual cu litera „ş” din limba turcă („s cu sedilă”). Diferența grafică dintre cele două semne diacritice este aproape insesizabilă pentru mărimi mici de text⁴ (vezi Figura 1).

Standardul Latin-2 nu a fost niciodată asociat limbii turce⁵. Cu toate acestea, caracterele „ș” și „ț” au fost definite în acest standard, pentru limba română, drept „s cu sedilă” (caracterul turcesc), respectiv „t cu sedilă” (caracter practic nefolosit în nici o limbă).⁶

¹ 9 MAI 2010, VERSIUNEA 2.0 (Cea mai actualizată versiune a acestui document, împreună cu alte resurse conexe, se vor găsi întotdeauna la adresa <http://www.moongate.ro/products/diacritice/>)

² Vezi http://www.secarica.ro/html/s-uri_si_t-uri.html.

³ Standardul ISO/IEC 8859-2, cunoscut și ca Latin-2.

⁴ Vezi și nota .

⁵ Caracterele pentru limba turcă au fost incluse inițial în standardul ISO/IEC 8859-3 (Latin-3), creat pentru Europa de Sud; câțiva ani mai târziu a fost creat un standard special pentru limba turcă, ISO/IEC 8859-9 (Latin-5).

⁶ Această discrepanță dintre uzanța autohtonă și standardele ISO a fost probabil cauzată de faptul că înseși documentele românești care descriau semnele diacritice respective le-au numit sedile timp de aproape două secole – vezi http://www.capisci.ro/articole/Sedile_%C3%AEn_rom%C3%A2n%C4%83

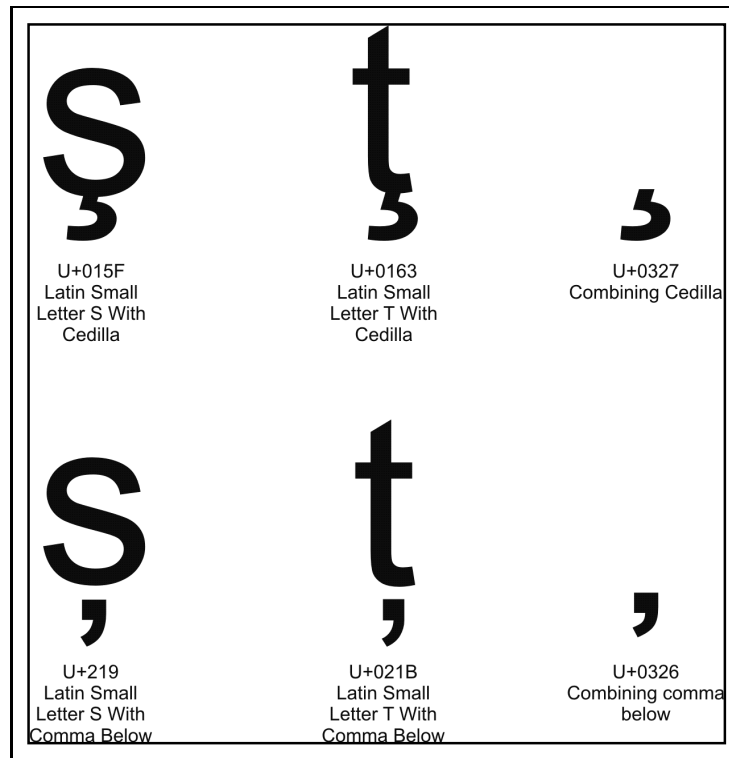


Figura 1: Caracterele în cauză, mărite pentru vizibilitatea semnelor diacritice. Primul rând conține caracterele „s cu sedilă”, „t cu sedilă” și semnul diacritic sedilă. Al doilea rând conține caracterele „s cu virgulă”, „t cu virgulă” și semnul diacritic virgulă.

În anul 1997 compania Apple a schimbat caracterele din standardul proprietar MacOS Romanian în așa fel încât să utilizeze semnele diacritice corecte pentru „ș” și „ț”.⁷ În același an Asociația de Standardizare din România a protestat pe lângă ISO în privința standardului Latin-2, însă singura modificare făcută un an mai târziu de către ISO a fost adăugarea unei note care permitea interpretarea semnelor respective drept „s cu virgulă”, respectiv „t cu virgulă”, iar asta numai în măsura în care expeditorul și destinatarul mesajului se puneau cumva de acord în această privință.⁸

În anul 1999 a apărut prima versiune a standardului Unicode care să conțină caracterele corecte românești „s cu virgulă” și „t cu virgulă”. Dificultățile de adoptare a standardului Unicode în formă completă au cauzat însă întârzieri de mai bine de zece ani în adoptarea sa pe scară largă.⁹

⁷ <http://unicode.org/Public/MAPPINGS/VENDORS/APPLE/ROMANIAN.TXT>

⁸ „Note - Subject to the agreement of originator and receiver, in information interchange the letters S and T WITH CEDILLA BELOW may be used to substitute for the letters S and T WITH COMMA BELOW”, http://www.secarica.ro/html/s-uri_si_t-uri.html

⁹ Caracterele cu semne diacritice corecte sunt incluse și în standardul ISO/IEC 8859-16 (Latin-10), publicat în 2001. Latin-10 este suportat în navigatoarele de Internet Firefox, Chrome și Opera și în sistemele de operare Mac OS, de la versiunea 10.4 (<http://www.opensource.apple.com/source/CF/CF-368.25/String.subproj/CFStringEncodingExt.h>), și Linux (<http://www.kernel.org/doc/man-pages/online/pages/man7/latin10.7.html>). Deși prezintă unele avantaje în fața Unicode (în special pentru că standardele ASCII extinse sunt single-byte), adoptarea pe scară largă a Unicode a făcut

Așa se face că timp de douăzeci de ani aproape toate textele scrise în limba română în medii informatice au fost scrise fie fără semne diacritice, fie cu semne diacritice greșite. Abia sistemul de operare Windows Vista, apărut la sfârșitul anului 2006, a fost primul sistem de operare utilizat pe scară largă de consumatorii de conținut care să folosească în mod nativ caracterele corecte pentru limba română. Chiar și așa, penetrarea lentă a acestui sistem de operare și a celor ulterioare face ca migrarea către utilizarea în practică a semnelor diacritice corecte să fie încă și astăzi, în 2010, mai mult un subiect de discuție sporadică decât obiectul unor acțiuni concrete.

2. Contextul de aplicabilitate al acestei lucrări

În această lucrare voi utiliza termenul de *colecție de text* în sens cât se poate de abstract și de cuprinzător: orice colecție de texte în limba română stocate electronic, indiferent de formatul concret de stocare sau de modul de prezentare. De la jurnal, revistă, carte sau enciclopedie tipărită până la mesaje e-mail, site-uri Internet sau etichete de text din cadrul aplicațiilor software, toate vor migra mai devreme sau mai târziu de la caractere cu sedile la caractere cu virgule.¹⁰

În ceea ce privește *factorii abstracți* care influențează alegerea momentului optim de migrare, am căutat să identific o structură suficient de generică încât să fie aplicabilă oricărei situații practice, în contextul definiției cuprinzătoare din paragraful anterior.

Pe de altă parte, *datele statistice concrete* prezentate în această lucrare sunt specifice numai colecțiilor de text care satisfac simultan următoarele criterii independente:

1. sunt consultate prin intermediul unui navigator de Internet (*web browser*);
2. sunt consultate de consumatori eterogeni în privința platformei software utilizate.

Dacă măcar unul dintre criteriile de deasupra nu se aplică, atunci trebuie *ignore* complet toate datele statistice concrete din această versiune a acestui document.¹¹ În acest caz trebuie utilizate resursele indicate la sfârșitul acestui document (dacă se aplică), sau trebuie căutate și adaptate datele statistice concrete asociate situației concrete la structura prezentată aici.

În privința dimensiunii temporale a deciziei de migrare am decis să nu includ niciun fel de date concrete, întrucât nivelul estimat de eroare al oricărei predicții de această natură ar fi prea mare pentru orice scop practic. Am ales în schimb să actualizez acest document și resursele conexe pe măsură ce evoluează situația (vezi nota **Error! Reference source not found.** sau ultimul paragraf din această lucrare).

ca, pentru moment cel puțin, Latin-10 să nu fie adoptat pe scară largă; un alt factor decisiv în această privință este și faptul că sistemele de operare și navigatoarele de la Microsoft nu suportă acest standard..

¹⁰ Eu personal am întâlnit această problemă în contextul discuțiilor de la Wikipedia în limba română; acolo mi-am și format și sintetizat în mare parte argumentele expuse în această lucrare, interacționând cu comunitatea de voluntari din cadrul proiectului respectiv.

¹¹ De exemplu dacă (1) este vorba despre o aplicație client-side, (2) este o aplicație online dedicată clienților care folosesc terminale mobile, (2) este o aplicație (online sau offline) care rulează numai pe o platformă anume, sau (1+2) este o aplicație client-side pentru terminale mobile.

3. Alegerea momentului: de ce este important

După cum am arătat mai sus, diferența grafică dintre cele două variante de caractere este în cea mai mare parte a timpului ne semnificativă.¹² Din acest motiv *nu există nicio presiune naturală considerabilă pentru adoptarea semnelor corecte* – practic toți consumatorii de conținut pot interpreta corect semnele diacritice „vechi”, iar majoritatea acestora nu sunt oricum la curent cu această problemă grafică minoră. Chiar și într-un sens mai profund chestiunea este la fel de neimportantă: *într-un text scris în limba română distincția dintre cele două tipuri de semne diacritice nu are valoare semantică*, deoarece utilizarea uneia dintre variante în defavoarea celeilalte nu aduce niciun plus de informație, indiferent de felul în care este interpretat textul.¹³

Prin urmare avem de-a face cu o *problemă semnificativă de interoperabilitate cauzată de rezolvarea unei probleme minore de prezentare*. Situația pare absurdă, însă faptul că nu există (și nu poate exista) nici o soluție alternativă pentru problema de prezentare legitimează problema de interoperabilitate.

În sistemele de operare ale companiei Microsoft anterioare Windows Vista caracterele cu diacritice corecte sunt vizibile numai în Windows XP, și asta numai în anumite condiții.¹⁴ Datorită cotei de piață uriașe a sistemelor de operare ale companiei Microsoft în rândul consumatorilor de conținut¹⁵, *aceste considerente fac migrarea către caracterele cu diacritice corecte o chestiune discutabilă în absența penetrării masive a sistemelor de operare Windows Vista sau mai noi pe piață*.

Pe de altă parte, unii factori de decizie ai diverselor colecții de text de limbă română vor fi în mod inevitabil *early adopters*¹⁶ ai noilor caractere cu diacritice corecte. Pe măsură ce trece timpul, pe măsură ce penetrează Windows Vista și sisteme de operare mai recente și pe măsură ce diverse colecții de text migrează la noile diacritice, va crește masa de conținut și de consumatori de conținut axați pe noile caractere. Odată ce se atinge o masă critică, *acei creatori de conținut care vor mai oferi text cu diacriticele incorecte vor fi văzuți ca depășiți*. Dacă în prezent există motive justificate pentru a amâna migrarea către diacriticele corecte¹⁷, odată ce se atinge masa critică *nu va mai exista nicio scuză pentru întârzierea migrării*: în ultimă instanță, diacriticele noi sunt cele corecte, iar cele vechi sunt pur și simplu incorecte în limba română!

Totuși voi arăta mai jos că independent de felul în care sunt văzuți din afară sau de corectitudinea tehnică a diacriticelor folosite în colecțiile lor de text, unii dintre creatorii de conținut de limbă română vor avea *motive întemeiate pentru a adopta noile diacritice înaintea celorlalți*, iar alții vor avea *motive întemeiate pentru a întârzia migrarea pentru o perioadă semnificativă* chiar și după momentul apariției masei critice. Scopul acestui

¹² Există totuși situații în care diferența este ușor de sesizat chiar și pentru un consumator nevizat, mai ales atunci când se folosesc mărimi mari de literă (titlul unei cărți pe copertă, titlurile de pe afișe sau materiale publicitare etc.)

¹³ Mai puțin cazul în care un text scris în română conține citate sau nume turcești. Chestiunea este discutată mai pe larg în secțiunea dedicată colecțiilor de text.

¹⁴ Chestiunea este analizată pe larg în secțiunea dedicată consumatorilor de conținut.

¹⁵ Vezi Tabelul 1

¹⁶ Persoane (fizice sau juridice) care doresc să adopte cât mai repede tehnologiile cele mai recente.

¹⁷ Pentru brevităte voi folosi în continuare sintagmele „diacritice corecte” și „diacritice noi” pentru „caractere care folosesc semnele diacritice corecte” (virgule), respectiv „diacritice incorecte” și „diacritice vechi” pentru „caractere care folosesc semnele diacritice vechi” (sedile).

FACTORII CARE INFLUENȚEAZĂ MOMENTUL OPTIM DE MIGRARE LA DIACRITICELE CORECTE ÎN LIMBA ROMÂNĂ

document este tocmai acela de a identifica factorii care influențează momentul în care trebuie luate aceste decizii, de la caz la caz.

4. Factori de influență

După cum am arătat în secțiunea , două forțe opuse acționează simultan asupra deciziei de migrare la diacriticele corecte:

- *Pentru migrare cât mai rapidă*: tehnologia este deja disponibilă, conținutul ar putea fi deja vizualizat de majoritatea consumatorilor, iar rezultatul ar fi utilizarea caracterelor corecte în limba română. În plus, imaginea ultimilor creatori de conținut care să migreze va avea probabil de suferit într-o oarecare măsură.
- *Pentru amânarea migrării*: diverse probleme de lizibilitate și interoperabilitate, dintre care unele foarte semnificative.

Aceste două forțe vor avea o evoluție dinamică de-a lungul timpului, în sensul că prima va crește în defavoarea celei de-a doua, până la eliminarea completă a acesteia din urmă.

După o analiză îndelungată a structurii diversilor factori care influențează aceste două forțe contrare, am identificat trei piloni pe care se sprijină întregul raționament:

- a) *Consumatorii de conținut*: cititorii, utilizatorii produselor software etc.
- b) *Creatorii de conținut*: edituri, deținători de site-uri, producători de software etc.
- c) *Colecțiile de text*: conținutul efectiv al revistelor, site-urilor, aplicațiilor etc.

În acest capitol voi analiza felul în care fiecare dintre acești trei piloni afectează fiecare dintre cele două forțe identificate mai sus.

4.1 Consumatorii de conținut

Este de la sine înțeles că cel mai important dintre cei trei piloni este cel reprezentat de consumatorii de conținut: indiferent de capacitățile tehnice ale creatorilor de conținut și de colecțiile lor de text, orice demers este inutil în măsura în care conținutul nu poate fi consumat sau, în cazul aplicațiilor interactive, consumatorul nu poate interacționa cu interfața în așa fel încât să obțină acces la conținutul propriu-zis.

Prima întrebare în ceea ce-l privește pe consumatorul de conținut este legată de modalitatea prin care acesta consumă în mod concret conținutul. *Dacă mediul final de consum nu implică tehnologii aflate sub controlul consumatorului, atunci consumatorul nu este un factor semnificativ* în luarea deciziei de migrare. În această situație se află conținutul prezentat exclusiv pe medii tipărite, sau în general pe orice medii în care consumatorul are un rol pasiv din punctul de vedere al tehnologiilor implicate (cărți, jurnale, reviste, afișe, prezentări video, filme și așa mai departe). În plus, conținutul consumat în condiții controlate de creatorii de conținut beneficiază în mare măsură de aceleași derogări (de exemplu aplicații care rulează în mod chioșc¹⁸, aplicații online sau client-side care rulează într-un mediu proprietar, conținut prezentat pe hardware dedicat precum cititoarele de cărți electronice).

¹⁸ Terminale dedicate, instalate în locuri publice, așa cum sunt cele din aeroporturi, gări, puncte de interes turistic etc.

Dacă însă mediul de consum este dependent de tehnologii controlate de consumator, atunci devin relevante două subcategorii de factori care influențează capacitatea consumatorului de a utiliza conținutul:¹⁹

1. *Lizibilitatea* – pot consuma conținutul?²⁰

2. *Interactivitatea* – pot interacționa cu interfața?

Pentru a cântări capacitatea de lizibilitate a consumatorilor în contextul diacriticelor corecte trebuie analizat nivelul de utilizare al platformelor care suportă diacriticele corecte în *contextul consumatorilor colecției de text* în speță.²¹

Datele statistice utilizate în acest document sunt următoarele:²²

Tabelul 1: Datele utilizate pentru generarea graficelor

Sistem de operare	Cota de piață	Afișează	Serie ²³
Windows XP ²⁴	73,14%	Lizibil (70,8%)	Simplu (34%)
		Ilizibil (29,2%)	Dificil (66%)
Windows Vista	8,91%	Perfect	Simplu
Windows 7	15,89%	Perfect	Simplu
Mac OS X	0,75%	Perfect	Simplu
Linux	0,74%	Nesigur	Simplu
Altele	0,73%	Nesigur	Dificil

Tabelul 2: Cota de piață a diverselor versiuni de Internet Explorer

Versiune Internet Explorer	Cota de piață
Internet Explorer 8	17,95%
Internet Explorer 7	9,75%
Internet Explorer 6	11,43%
Altele	0,11%

¹⁹ Cel mai reprezentativ exemplu în acest sens sunt colecțiile de text ale site-urilor și aplicațiilor Internet/intranet. În aceeași situație se află însă orice colecție de text care poate fi interpretată pe mai multe platforme software – aplicații, documente distribuite, mesaje e-mail (e.g. newsletters) și așa mai departe.

²⁰ Includ în această subcategorie și problemele de accesibilitate pentru persoanele cu dizabilități – de exemplu în ce măsură aplicațiile care citesc textul pentru consumatorii nevăzători sunt capabile să recunoască și să interpreteze corect texte scrise cu diacriticele noi.

²¹ Această precizare este crucială pentru o analiză corectă în cazul unor situații particulare; vezi și nota .

²² Datele din Tabelul 1 și Tabelul 2 sunt colectate în mai 2010 de la <http://gs.statcounter.com/>, pentru consumatorii din România. Pentru o analiză riguros exactă a consumatorilor de conținut în limba română la nivel global ar fi necesare datele statistice asociate acestui grup demografic specific, independent de poziția geografică (consumatorii de conținut în limba română nu se găsesc numai în România, ci și în Republica Moldova și pe alte meridiane; în plus, nu toți consumatorii de conținut din România preferă limba română). Totuși, așa după cum am indicat în mod repetat în acest document, fiecare creator de conținut trebuie să țină cont de capabilitățile tehnice ale propriilor consumatori. Prin urmare am considerat acceptabile aproximările rezultate din utilizarea unui criteriu geografic pentru ilustrarea orientativă a situației curente în acest document.

²³ Statistică relevantă în special pentru secțiunea legată de creatorii de conținut.

²⁴ Aproximări semnificative în ambele subcategorii. Frațiile din subcategorii sunt procente din cota Windows XP.

FACTORII CARE INFLUENȚEAZĂ MOMENTUL OPTIM DE MIGRARE LA DIACRITICELE CORECTE ÎN LIMBA ROMÂNĂ

În privința Windows XP, care deocamdată rămâne lider detașat (nu numai că este cea mai utilizată platformă, dar este mai utilizată decât toate celelalte la un loc), situația este la fel de incertă, însă merită investigată. La instalare, Windows XP nu suportă deloc diacriticele corecte – pur și simplu pe ecran apar niște „pătrățele” în locul caracterelor respective (ultima coloană din Tabelul 3). Există două modalități principale de a obține compatibilitate cu diacriticele corecte în Windows XP²⁵:

- Prin instalarea explicită a unui pachet software suplimentar de la Microsoft²⁶.
- Prin instalarea Internet Explorer 7 sau Internet Explorer 8, aplicații care sunt în mod normal instalate prin mecanismul de actualizate automată al Windows.

Pentru prima opțiune nu avem la dispoziție statistici, însă este destul de puțin probabil că un consumator oarecare de conținut a făcut efortul să descarce și să instaleze un astfel de pachet software. Pe de altă parte, actualizarea navigatorului Internet Explorer este una automată (și dezirabilă pentru motive independente de diacritice), deci este de așteptat ca o parte semnificativă a consumatorilor să fi instalat această actualizare.

Figura 2: Gradul de utilizare al diverselor sisteme de operare

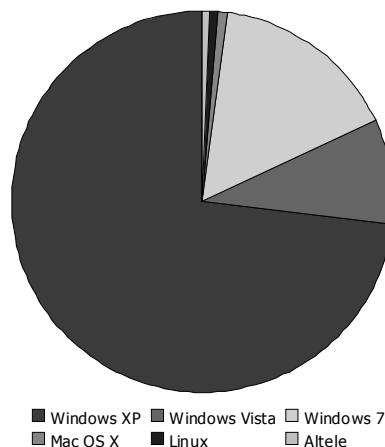
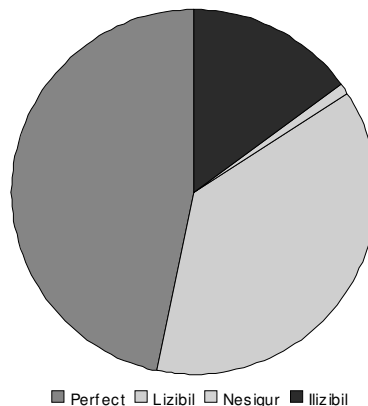


Figura 3: Capacitatea sistemelor de operare de a afișa diacriticele



²⁵ <http://www.microsoft.com/Romania/Diacritice.aspx> – în realitate atât Internet Explorer 8 cât și Internet Explorer 7 sunt capabile să opereze substituția de corp de literă pentru a obține rezultatele de pe coloana a doua din Tabelul 3.

²⁶ „Actualizare de fonturi corespunzătoare extinderii Uniunii Europene” de la <http://www.microsoft.com/downloads/details.aspx?FamilyID=0ec6f335-c3de-44c5-a13d-a1e7cea5ddea&DisplayLang=ro>

Totuși asta înseamnă că mai bine de o treime dintre utilizatorii de Internet Explorer nu pot vedea text scris cu diacriticele corecte, atâta timp cât utilizează Windows XP. Nu avem la dispoziție date statistice globale de încredere care să coroboreze sistemul de operare cu navigatorul utilizat²⁷, deci vom fi nevoiți să presupunem că toți utilizatorii de Windows XP se supun aceleiași proporții identificate pentru Internet Explorer, indiferent dacă utilizează acest navigator sau altul.

Deja în 2010 *majoritatea consumatorilor de conținut pot să citească text care folosește diacriticele corecte*. Există totuși câteva rezerve semnificative:

- Utilizatorii pentru care textul este doar lizibil (nu perfect) pot citi textul, însă, în funcție de situație, *unii vor observa o diferență sesizabilă, inestetică și deranjantă de afișare a caracterelor respective* (a doua coloană din Tabelul 3). Suportul pentru diacriticele corecte se limitează în Windows XP la numai câteva corpuri de literă – pentru celelalte, sistemul de operare substituie pur și simplu caracterele respective cu aceleași caractere din cele mai apropiate corpuri de literă pentru care dispune de caracterele în speță. Rezultatul este cel așteptat: textul este lizibil, însă în funcție de diferența vizuală dintre corpul de literă utilizat în text și cel disponibil rezultatele pot fi inestetice (în tabel am folosit Arial Black).
- Chiar în cazul colecțiilor de text accesibile via Internet este posibil ca unele aplicații specifice să se adreseze în mod particular consumatorilor care utilizează o paletă relativ îngustă de sisteme de operare. Un exemplu relevant sunt aplicațiile sau subdomeniile dedicate pentru platforme mobile.²⁸

	Windows 7, Vista	Windows XP (nou)	Windows XP (vechi)
Diacritice corecte	arșiță	arșiță	ar i i ă
Diacritice vechi	arșiță	arșiță	arșiță

Tabelul 3: Afișarea celor două tipuri de diacritice în funcție de gradul de lizibilitate al diacriticelor noi (de notat că diacriticele vechi sunt perfect lizibile indiferent de context)

²⁷ De fapt statistica ideală ar fi chiar cea pe care încerc să o estimez aici (lizibilitate perfectă/lizibilitate limitată/ilizibilitate); în lipsa ei, cea mai bună aproximare ar fi o coroborare a lizibilității limitate (gradul de instalare al EUupdate.EXE sau IE7 sau IE8) cu sistemul de operare Windows XP. Totuși procentele sunt deocamdată suficient de mari în toate categoriile încât aproximările din text să nu afecteze în mod semnificativ calitatea analizei.

²⁸ Multe site-uri publice oferă alternative pentru platforme mobile. De pildă un consumator care accesează pentru prima dată pagini din domeniul <http://www.moongate.ro/> folosind un dispozitiv mobil este redirecționat automat către pagina corespunzătoare din domeniul <http://m.moongate.ro/>; în cazul acestui al doilea domeniu sunt relevante numai statisticile legate de platformele mobile.

FACTORII CARE INFLUENȚEAZĂ MOMENTUL OPTIM DE MIGRARE LA DIACRITICELE CORECTE ÎN LIMBA ROMÂNĂ

Acestea fiind spuse, trebuie menționat în mod proeminent că indiferent de platforma specifică a consumatorilor, indiferent de publicul țintă al colecției de text și indiferent de numărul lor, *fracțiunea consumatorilor care nu pot vizualiza diacriticele corecte se vor afla practic în imposibilitate de a consuma conținutul*. Iar alternativa este afișarea aceluiași text, utilizând caractere cu semne diacritice tehnic incorecte, dar care sunt aproape identice din punct de vedere vizual și pe care le poate citi oricine (vezi al doilea rând din Tabelul 3). În ultimă instanță trebuie pusă în balanță corectitudinea academică față de pierderea de facto a unei fracțiuni a cititorilor.

Celălalt factor care trebuie luat în considerare este capacitatea consumatorilor de a interacționa cu interfața colecției de text. Cea mai proeminentă funcție a interfeței în această privință este funcționalitatea de căutare: dacă colecția de text conține diacritice corecte iar consumatorul operează o căutare utilizând diacriticele vechi (sau viceversa)²⁹ atunci consumatorul nu va obține rezultatele dorite. *Practic toate problemele de interacțiune pot fi rezolvate prin soluții tehnice*, însă acestea trebuie luate în considerare și rezolvate din timp.³⁰ *Totuși problemele de interacțiune sunt printre puținele care pot fi rezolvate înainte de începerea migrării și ar trebui rezolvate cât mai curând.*³¹

4.2 Creatorii de conținut

Al doilea pilon de influență al deciziei de migrare este capacitatea creatorilor de conținut de a crea conținut folosind diacriticele corecte. *Este prea puțin important dacă cititorii pot citi, atâta vreme cât scriitorii nu pot scrie.*

Prin urmare factorii care influențează creatorii de conținut sunt dictați în primul rând de o logică similară celei legate de consumatori:

Poate autoritatea sub egida căreia se generează conținut să controleze mijloacele tehnice ale creatorilor individuali de conținut?

Răspunsurile posibile la această întrebare sunt cu mult mai variate decât în cazul consumatorilor. Am ales aici numai cazurile extreme, pentru exemplificare:

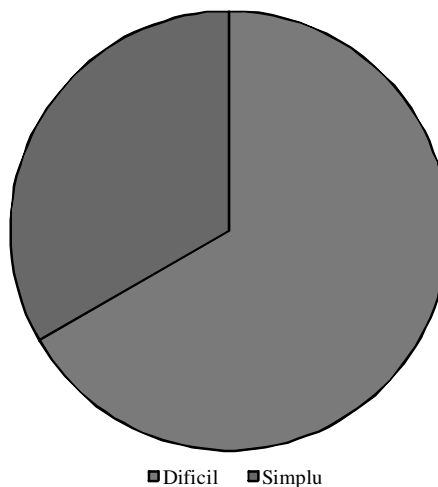


Figura 4: În ce măsura se poate scrie text cu diacriticele corecte

²⁹ Vezi și ultima coloană din Tabelul 1, reprezentată grafic în Figura 3.

³⁰ Aproape orice problemă de acest fel poate fi rezolvată prin interpretarea flexibilă a datelor de intrare, așa cum procedează Google sau DEX online.

³¹ O colecție de text proeminentă de limbă română care a migrat foarte curând la diacriticele corecte este DEX online (<http://dexonline.ro/>). Au fost însă luate în calcul aproape toate problemele identificate în acest document: consumatorii care nu pot citi diacriticele corecte pot alege să consulte dicționarul folosind diacriticele vechi, iar interacțiunile cu colecția de text ignoră în mod implicit semnele diacritice din textul de intrare. Singura scăpare este legată de limba turcă, în sensul că și cuvintele care ar trebui scrise cu sedilă au fost convertite automat la forma cu virgulă (e.g. <http://dexonline.ro/definitie/siret> versus <http://tr.wiktionary.org/wiki/%C5%9Ferit>). Totuși absența unei necesități de interoperabilitate ulterioară internă dintre dicționar și orice altă colecție de text face ca această scăpare să fie lipsită de orice efecte adverse concrete.

- *Redactori tradiționali* (jurnal, revistă, carte tipărită ș.a.m.d.): acești redactori lucrează la sediul societății care i-a angajat, iar societatea respectivă are control complet asupra platformei software utilizate la sediu. Chiar și redactorii care preferă să folosească mijloace tehnice proprii pot fi constrânși să urmeze standardele societății (e.g. atunci când redactează text utilizând laptopul sau computerul propriu). În acest caz creatorii de conținut sunt un factor neglijabil, deoarece autoritatea decizională va opera migrarea pe baza celorlalți doi piloni.
- *Redactori voluntari, independenți* (de exemplu redactori la Wikipedia sau alte proiecte similare, persoanele care contribuie cu comentarii în diverse site-uri, site-uri de socializare ș.a.m.d.): acesta este unul dintre cei mai puternici factori pentru amânarea adoptării diacriticele corecte (vezi dedesubt).

După cum se vede în Figura 4, *mai bine de jumătate dintre creatorii independenți de conținut au dificultăți în a scrie text utilizând diacriticele corecte*.³² Motivul central pentru această limitare este faptul că utilizatorii de Windows XP, deocamdată majoritari în statisticile globale, nu pot genera conținut utilizând diacriticele corecte decât în condiții destul de stricte.³³ În aceste condiții decizia trebuie luată în concordanță cu capacitatea factorilor decizionali de a influența mijloacele tehnice utilizate de creatorii individuali de conținut, sau de a implementa soluții tehnice pentru ameliorarea sau rezolvarea acestei probleme specifice. În cazul proiectelor bazate pe voluntariat, așa cum este Wikipedia, este evident că migrarea în condițiile actuale nu ar avea sorți de izbândă, chiar lăsând la o parte ceilalți piloni, în absența unor soluții tehnice specifice.³⁴

4.3 Colecția de text

Însăși colecția de text, cel mai puțin important dintre cei trei piloni identificați aici, va fi pentru mulți creatori de conținut impedimentul major în decizia de migrare, chiar și atunci când migrarea ar fi dezirabilă din celelalte puncte de vedere. Orice analiză a deciziei de migrare către diacriticele corecte este în mod firesc axată în primul rând pe utilizabilitate din punctul de vedere al consumatorului de conținut, în al doilea rând pe capacitatea creatorului de conținut de a genera conținut și abia în ultimul rând pe *disponibilitatea creatorului de conținut de a migra în mod retroactiv conținutul existent la diacriticele corecte*.

Colecția de text influențează decizia de migrare în funcție de următorii factori:

1. *Relevanță*: este relevant să vă puneți întrebări legate de migrarea colecției de text numai în măsura în care aceasta este vizibilă. De exemplu arhiva digitală a unui periodic publicat exclusiv în formă tipărită nu face în general obiectul migrării. Evident, dacă există o parte accesibilă în format electronic și una stocată intern atunci numai partea accesibilă este relevantă.
2. *Mărime*: semnificația acestui factor trebuie coroborată în general cu următorul, complexitatea colecției de text. Există însă un caz particular în care contează exclusiv mărimea: atunci când ea este nulă. Dacă un creator de conținut începe

³² Datele sunt extrase de pe ultima coloană din Tabelul 1.

³³ <http://www.stefamedia.ro/diacritice-romanesti-corecte-in-windows-xp/>

³⁴ La Wikipedia s-a luat deja decizia migrării la diacriticele corecte, deoarece au fost identificate soluții tehnice concrete care permit vastei majorități a redactorilor individuali să genereze conținut utilizând diacriticele corecte indiferent de platforma software pe care o utilizează. Pentru detalii despre acest caz particular vezi secțiunea dedicată studiului de caz Wikipedia.

FACTORII CARE INFLUENȚEAZĂ MOMENTUL OPTIM DE MIGRARE LA DIACRITICELE CORECTE ÎN LIMBA ROMÂNĂ

lucrul la o colecție de text complet nouă atunci trebuie cântărită cu atenție opțiunea de a utiliza de la bun început diacriticele corecte – în măsura în care aceasta este o opțiune acceptabilă din celelalte puncte de vedere atunci ar trebui adoptată, întrucât astfel se va evita mai târziu costisitorul proces de migrare retroactivă.

3. *Complexitate*: am văzut mai sus că distincția dintre caracterele cu virgulă și cele cu sedilă nu are valoare semantică în limba română. Altfel spus, un text scris în exclusivitate în limba română ar putea fi convertit cu ușurință la diacriticele corecte printr-o simplă operațiune de căutare și înlocuire automată. În general acest lucru este adevărat, însă cu unele rezerve pe care le voi analiza dedesubt.

Complexitatea colecției de text este o măsură a frecvenței situațiilor care necesită intervenție umană. În cazul colecției de text de la Wikipedia am identificat câteva tipologii specifice de situații problematice, probabil reprezentative pentru cazul general:

- *Texte scrise în altă limbă* (cel mai notabil în turcă) fără notație adecvată.³⁵ Dacă textele (inclusiv numele de persoane, locuri, evenimente ș.a.m.d.) în turcă sunt marcate explicit ca fiind scrise în turcă există posibilitatea de a automatiza procesul prin evitarea schimbării semnelor diacritice în cazul acestor fragmente de text.
- *Identificatori de resurse care conțin semne diacritice*. Pe lângă cazul general (URI) mai pot exista o sumedenie de identificatori interni sau externi a căror integritate structurală trebuie menținută de-a lungul procesului de migrare, în ambele sensuri.³⁶
- *Interoperabilitatea cu alte colecții de text*, în special în condițiile unei migrări parțiale.³⁷

5. Studiu de caz: Wikipedia în limba română

Am fost implicat în aproape toate discuțiile de la Wikipedia în limba română în privința deciziei de migrare la diacriticele corecte, iar mărimea și complexitatea colecției de text, numărul mare și eterogenitatea redactorilor și consumatorilor de conținut ai acestui proiect îl fac probabil unul dintre cele mai bune studii de caz posibile în contextul acestui document.

La Wikipedia în limba română, discuțiile despre migrarea la diacriticele corecte au început încă din anul 2007. Deși la acel moment au fost considerate premature, în urma discuției a fost creată o pagină dedicată subiectului.³⁸ La sfârșitul aceluiași an a fost

³⁵ De exemplu în HTML: <http://www.w3.org/TR/WCAG10-HTML-TECHS/#language> [en]

³⁶ Cele mai la îndemână exemple sunt legate de colecția de text de la Wikipedia în limba română. Printre identificatorii interni de resurse se numără legăturile interne între articole și cele care leagă articole despre același subiect în mai multe limbi. Identificatorii externi conținuți în Wikipedia sunt legăturile (URI) către pagini de pe alte site-uri și care pot conține caractere cu diacritice. Identificatorii externi pe care trebuie să-i gestioneze Wikipedia sunt legăturile (URI) dinspre alte site-uri către articolele din Wikipedia care pot conține caractere cu diacritice; aceasta este mai degrabă o responsabilitate morală datorată numărului relativ mare de documente care fac trimitere la enciclopedie.

³⁷ De exemplu atunci când colecția de text are relevanță parțială în privința diacriticelor, dar partea publică a colecției de text (cea care urmează să fie migrată) interacționează prin sisteme automate cu partea istorică/privată/confidențială care nu este migrată.

³⁸ http://ro.wikipedia.org/wiki/Wikipedia:Corectarea_diacriticelor

implementată o soluție tehnică de natură să forțeze utilizarea diacriticelor *vechi*, în scopul uniformizării conținutului (unii redactori deja scriau conținut folosind diacriticele noi). La sfârșitul lui 2008 și începutul lui 2009 discuțiile au început să se orienteze către identificarea problemelor și soluțiilor tehnice concrete asociate utilizării diacriticelor corecte și s-au început câteva demersuri tehnice concrete, deși fără efecte pentru consumatori, în vederea viitoarei migrări.

Dată fiind absența unor statistici și analize concrete, comunitatea a decis la 1 martie 2010 să implementeze un proiect pilot prin care se permitea utilizarea diacriticelor corecte pe cele mai multe dintre paginile interne ale proiectului, cu scopul de a identifica numărul de redactori incapabili să le citească; nu s-a înregistrat nicio plângere.

Pe fondul discuțiilor de la Wikipedia și în urma invitației organizatorilor ConsILR de a participa la conferință, am organizat în perioada 27 aprilie–4 mai 2010 un sondaj național³⁹ în scopul identificării capabilităților tehnice ale consumatorilor de limbă română în privința lizibilității diacriticelor noi. Discuțiile din timpul sondajului au dus la identificarea unui criteriu concret pe baza căruia urma să se ia decizia migrării la diacriticele corecte.⁴⁰

Sondajul a beneficiat de un nivel adecvat de expunere și participare⁴¹, dar am determinat ulterior că voturile exprimate au fost eronate în proporție mult prea mare pentru orice scop practic.⁴² Totuși datele statistice obținute prin analizarea vizitatorilor (ignorând voturile exprimate) au confirmat faptul că distribuția capabilităților tehnice ale vizitatorilor Wikipedia în limbă română este foarte apropiată de datele statistice naționale prezentate în acest document. Coroborând această informație cu soluțiile tehnice identificate anterior, de natură să amelioreze problema lizibilității pentru segmentul consumatorilor afectați, am constatat că a fost satisfăcut criteriul convenit în timpul desfășurării sondajului. Prin urmare, pe 5 mai 2010 s-a luat decizia de migrare cât mai rapidă la diacriticele corecte.⁴³

În cazul colecției de text de la Wikipedia în limba română procesul va fi cu siguranță unul de durată, însă analiza și execuția până în acest moment au urmat o traiectorie foarte apropiată de cea ideală. La Wikipedia, calitatea analizei și a deciziilor în această privință se datorează în cea mai mare măsură faptului că deciziile se iau prin consens, iar viziunile divergente ale diverșilor participanți la discuție au asigurat în cadrul discuțiilor o reprezentare optimă a celor două forțe identificate în secțiunea

6. Concluzii

Toate colecțiile de text care conțin text în limba română vor urma inevitabil, mai devreme sau mai târziu, standardul corect în privința semnelor diacritice. Singurul punct

³⁹ <http://www.moongate.ro/products/diacritice/sondaj/>

⁴⁰ „Mai puțin de 3% de vizitatori pe care nu îi putem ajuta (nu văd diacritice sau văd majuscule și au alt SO decât Windows sau Windows mai vechi de XP)”, la pagina menționată în nota .

⁴¹ <http://www.moongate.ro/products/diacritice/sondaj/date.php#statistici>

⁴² Aproximativ 20% dintre respondenții care au votat că nu pot citi diacriticele noi utilizau Windows Vista sau Windows 7, sisteme de operare despre care știm cu siguranță că în realitate le afișează perfect.

⁴³ http://ro.wikipedia.org/w/index.php?title=Discu%C5%A3ie_Wikipedia:Sfatul_B%C4%83tr%C3%A2nilor&oldid=3918861#Concluzii_finale

FACTORII CARE INFLUENȚEAZĂ MOMENTUL OPTIM DE MIGRARE LA DIACRITICELE CORECTE ÎN LIMBA ROMÂNĂ

delicat este alegerea momentului optim pentru migrare. Diferența dintre cele două variante este mică din punct de vedere vizual, dar există dificultăți tehnice potențiale semnificative asociate migrării.

Am căutat să identific aici factorii care influențează alegerea pragmatică a momentului optim de migrare pe baza următoarei structuri:

- Consumatorii de conținut
 - în ce măsură tehnologia se află sub controlul consumatorului
 - în ce măsură pot consuma conținutul
 - în ce măsură pot interacționa cu interfața
- Creatorii de conținut
 - în ce măsură pot controla tehnologia utilizată de redactorii individuali
 - în ce măsură redactorii individuali pot crea text folosind diacriticele corecte
- Caracteristicile colecției de text
 - relevanța
 - mărimea
 - complexitatea

Creatorii de conținut de limbă română ar trebui să ia cât mai curând următoarele măsuri concrete:

- Implementarea măsurilor tehnice necesare pentru a permite consumatorilor care folosesc deja diacriticele corecte să interacționeze cu interfața (e.g. pentru funcțiile de căutare).
- Implementarea măsurilor tehnice necesare în vederea migrării, dacă este cazul (în particular identificarea explicită a textelor scrise în limba turcă).
- Analiza situației propriului caz particular, prin prisma acestui document.
- Identificarea factorilor critici aplicabili colecției de text analizate.
- Stabilirea unui criteriu concret pentru migrare, pe baza factorilor aplicabili.
- Planificarea migrării (metodologie tehnică, estimare buget și resurse necesare).
- Monitorizarea și previzionarea evoluției situației în raport cu criteriul ales, în așa fel încât să poată alocă resursele necesare concretizării migrării în timp util.

Ultima versiune a acestui document, precum și alte noutăți și resurse suplimentare în materie se găsesc la adresa <http://www.moongate.ro/products/diacritice/>

Mulțumiri organizatorilor și referenților *ConsILR*, pentru că mi-au oferit imboldul inițial și încurajarea ulterioară pentru cristalizarea acestui document în formă scrisă; domnului *Cristian Secară*, pentru informațiile concrete care m-au ajutat să corectez unele statistici care altfel ar fi fost cu siguranță greșite; domnului *Cristian Adam*, pentru notele legate de Latin-10; *colectivității Wikipedia* în limba română, pentru discuțiile întotdeauna deschise; în acest context le mulțumesc în particular redactorilor *Cezarika1* (primul care a ridicat această problemă și ne-a arătat de ce merită discutată), *Strainu* (un early adopter prin definiție, foarte capabil pe partea tehnică, cel care a împins constant în direcția adoptării cât mai rapide a diacriticelor corecte la Wikipedia și care m-a ajutat și cu sugestii în privința acestui document), *Danutz* (care mi-a indicat site-ul

<http://gs.statcounter.com/>, utilizat pentru toate statisticile din acest document; în versiuni mai vechi ale acestui document am folosit statistici globale, semnificativ diferite), *AdiJapan* (ca întotdeauna, a reușit să găsească un echilibru între interlocutorii mai tehnici și cei mai puțin tehnici, între cei avangardiști și cei conservatori); și, independent de Wikipedia, *tatălui meu și prietenilor* care m-au ajutat să dau forma curentă acestei lucrări (știți voi cine sunteți).

CONSTRUCȚIA AUTOMATĂ DE CORPUSURI MULTILINGVE

TIBERIU BOROȘ, DAN TUFIȘ, ALEXANDRU CEAUȘU

Institutul pentru Cercetarea Inteligenței Artificiale – Academia Română

{tibi,tufis,aceausu}@racai.ro

Rezumat

Insuficiența resurselor lingvistice pentru multe dintre limbile naturale este principalul impediment în progresul tehnologic al prelucrării automate a acestor limbi. Succesul abordărilor statistice în traducerea automată pentru limbile de largă circulație, bine echipate sub raport cantitativ și calitativ cu resurse lingvistice, a evidențiat încă o dată necesitatea colectării și preprocesării corespunzătoare a unor volume cât mai mari posibile de resurse lingvistice. Pentru aplicațiile multi- și cros-linguale corpusurile paralele și (mai recent) cele comparabile sunt resurse primare indispensabile. Articolul de față prezintă o serie de unelte ce pot fi folosite în extragerea automată de corpusuri paralele sau puternic comparabile.

1. Introducere

În ultimii 15-20 de ani abordările bazate pe tehnici inductive și volume mari de date, ca și creșterea spectaculoasă a performanțelor de calcul și de stocare ale noilor generații de calculatoare au condus la progrese greu de anticipat în anii de început ai lingvisticii computaționale și ai traducerii automate. Utilizarea web-ului ca resursă primară de date, fără nici un fel de prelucrare lingvistică majoră, a demonstrat că utilizarea tehnologiilor statistice aplicate unor volume foarte mari de texte poate fi un răspuns la marea provocare a traducerii automate între toate limbile documentelor din spațiul web. Google translate¹ deja oferă traducere automată pentru 52 de limbi (2652 de perechi limbă sursă-limbă țintă). Sistemul Bing Translator² de la Microsoft oferă traducere automată pentru 30 de limbi (870 de perechi limbă sursă-limbă țintă). Pe lângă aceste exemple de sisteme publice foarte cunoscute pot fi amintite și companii mai noi³, focalizate pe traducere automată care dezvoltă sisteme comerciale, asigurând însă pentru unele perechi de limbi o calitate a traducerii aproape de nivelul traducerii umane⁴.

Dacă metodele de prelucrare numerică bazate exclusiv pe texte neprelucrate lingvistic (în engleză acest tip de prelucrare este numit *number crunching NC*) sunt aplicabile pentru orice pereche de limbi, în schimb volumul datelor necesare pentru o calitate acceptabilă a traducerii este uriaș (pentru unele perechi de limbi nici măcar conținutul actual al web-ului nu este suficient), necesitând resurse de calcul și stocare ce depășesc dotările grupurilor obișnuite de cercetare. Numai marile companii își pot permite dezvoltări de sisteme de prelucrare multilingve pentru un număr mare de perechi

¹ <http://translate.google.com/#> (consultare la data de 10 aprilie 2010)

² <http://www.microsofttranslator.com/Default.aspx> (consultare la data de 10 aprilie 2010)

³ http://www.translationguide.com/translation_company_links.php (o listă parțială la data de 10 aprilie 2010)

⁴ <http://www.languageweaver.com/page/home/>

arbitrare de limbi. Pe de altă parte cercetări recente (Koehn et al., 2007), (Hoang et al., 2009), (Tufiș et al., 2009) ș.a. au arătat că atunci când textele sunt prelucrate lingvistic (segmentate, dezambiguizate morpho-lexical, lematizate, parsate) cu un volum de date mult mai mic decât conținutul (pentru o anumită pereche de limbi) întregului web se pot obține traduceri comparabile și chiar superioare metodelor de tip NC. Fără îndoială însă că și pentru astfel de abordări cu cât volumul datelor primare este mai mare calitatea traducerilor crește.

Pentru sistemele de traducere automată bazate pe metode statistice (indiferent dacă sunt de tip NC sau includ prelucrări lingvistice) cele mai valoroase resurse textuale sunt corpusurile paralele ce conțin texte în două sau mai multe limbi astfel încât ele reprezintă traduceri reciproce. Unul dintre primele corpusuri paralele bilingv (engleză-franceză) este corpusul Hansards conținând transcrieri ale dezbaterilor în Parlamentul Canadei în perioada 1975-1988 (Germann, 2001). El a constituit întotdeauna o resursă fundamentală pentru studiile cros-linguale pentru perechea de limbi engleză-franceză.

Multilingvismul este o caracteristică esențială a Europei și implicațiile culturale, sociale și economice au ridicat multilingvitatea la nivelul unei preocupări politice de prim rang și absolut toate instituțiile Uniunii Europene produc un volum foarte mare de documente paralele. Necesitatea accesului cercetătorilor din domeniul prelucrării limbajelor naturale la aceste documente a determinat factorii responsabili (de exemplu OPOCE – biroul pentru publicații oficiale ale Uniunii Europene) să facă publice formatul electronic al unui volum din ce în ce mai mare de documente paralele pe baza cărora cercetătorii au construit primele corpusuri paralele europene publice. Primul dintre acestea a fost EuroParl (Koehn, 2005), disponibil pentru 11 din limbile comunității europene încă din 2001. El conține transcrieri ale sesiunilor din Parlamentul European din perioada 1996-2001. Versiunile următoare⁵ au extins cantitatea de texte, dar nu și numărul de limbi. Câțiva ani mai târziu, în 2006, a fost creat și distribuit prima variantă a corpusului JRC-Acquis (Steinberger et al., 2006). Versiunea V3 este în prezent⁶ cel mai mare corpus multilingv disponibil, conținând texte de natură juridică în 22 de limbi ale Uniunii Europene. Acest corpus conține peste 1 miliard de cuvinte, în medie 48 milioane de cuvinte/limbă. Din perspectiva limbii române, alături de JRC-Acquis mai sunt relevante corpusurile EMEA și OpenSubs (Tiedemann, 2009). EMEA conține documente ale Agenției Europene a Medicamentului în 22 de limbi, iar OpenSubs conține subtitrări în 30 de limbi. Textele românești din JRC-Acquis și EMEA conțin diacritice în schimb ele lipsesc din OpenSubs ceea ce face acest ultim corpus mai puțin util pentru prelucrarea limbii române.

La o privire generală, cele mai multe corpusuri paralele au două mari probleme inerente: conțin, în principal, limbi de largă circulație și sunt specifice unui anumit domeniu.

Folosirea de corpusuri comparabile devine o alternativă viabilă pentru antrenarea sistemelor de traducere în cazul domeniilor și limbilor care au o mică reprezentare în corpusuri paralele. Un corpus comparabil reprezintă o colecție de texte în două sau mai

⁵ Versiunea v5, distribuită la începutul anului 2010, acoperă perioada 1996-2009 și conține aproximativ 55 milioane de cuvinte pentru fiecare din cele 11 limbi: franceză, italiană, spaniolă, portugheză, engleză, olandeză, germană, daneză, suedeză, greacă și finlandeză.

⁶ <http://langtech.jrc.it/JRC-Acquis.html>

multe limbi care deși nu reprezintă traduceri reciproce riguroase, conțin informații similare. Aceste tipuri de corpusuri au grade de comparabilitate diferite (Skandina et al., 2010), cele mai utile fiind cele clasificate drept „puternic comparabile” (eng. *strongly comparable*). Aceste corpusuri comparabile conțin informații despre aceleași lucruri, folosesc același registru lingvistic și au un grad ridicat de suprapunere la nivelul echivalențelor lexicale de traducere. Deși antrenarea sistemelor de traducere folosind corpusuri comparabile implică un nivel de complexitate ridicat față de antrenarea tradițională ce folosește resurse paralele, textele comparabile sunt disponibile în proporție mult mai mare decât cele paralele.

2. Combinarea modelelor de traducere extrase din corpusuri paralele și corpusuri comparabile

În cadrul proiectului național STAR (PNII – IDEI 742/19.01.2009) la Institutul de Cercetări pentru Inteligența Artificială a fost realizat prototipul unui sistem de traducere pentru perechea de limbi engleză-română (Ceașu, 2009) folosind platforma Moses (Koehn et al, 2007) și serviciile de procesare a textului (Tufiș et al, 2007). Sistemul a fost antrenat pe un corpus de aproape un milion de unități de traducere cu peste treizeci de milioane de atomi lexicali. Corpusul este compus în proporție de 75% din texte din domeniul juridic, 5% texte din domeniul jurnalistic. Restul de 20% sunt echivalenți de traducere și unități de traducere extrase din ontologia lexicală românească Ro-Wordnet extinsă cu terminologie juridică (Tufiș et al, 2008). Experimentele noastre au arătat că traducerea textelor din domeniul juridic are o calitate bună dar, deși inteligibile, traducerile textelor din alte domenii (de ex. sport, medicină, turism etc.) au o calitate mult mai slabă. O soluție rațională pentru remedierea acestui aspect constă în construcția mai multor modele statistice de limbă și de traducere, fiecare caracteristic unui anumit registru textual. În continuare, un modul preliminar ar putea clasifica un text nou ce urmează a fi tradus într-una din clasele pentru care există modele de traducere și sistemul va efectua traducerea folosind modelul specific. O altă variantă constă în a combina diferitele modele, atribuind dinamic, în funcție de domeniul textului ce urmează a fi tradus, ponderi de influență diferite modelelor statistice combinate.

Proiectul european Accurat⁷ (248347/FP7), lansat la începutul acestui an, are ca obiectiv, printre altele, construcția unor corpusuri comparabile cât mai mari pentru care să poată complementa puținele corpusuri paralele disponibile pentru limbile proiectului (inclusiv româna). Combinarea modelelor de traducere construite din cele două categorii de corpusuri se va face studiind variantele amintite anterior și alegând soluția cea mai performantă.

Pentru a extinde sistemul de traducere și pentru alt domeniu decât cel juridic, decodorul Moses îi pot fi adăugate modele de traducere antrenate folosind corpusurile comparabile.

Moses este un decodor pentru sistemele de traducere automată care extinde traducerea formei de ocurență a cuvintelor cu modele de traducere factorizate. Spre deosebire de decodarea bazată pe echivalenții de traducere constituiți din secvențe contigue de cuvinte, care se bazează numai pe forma de ocurență a cuvintelor, traducerea factorizată

⁷ <http://www accurat-project.eu/>

poate lua în considerare informații suplimentare asociate atomului lexical, cum ar fi partea de vorbire, forma de dicționar a cuvântului sau descrierea sa morfo-sintactică.

În conformitate cu modelul traducerii factorizate, procesul de traducere presupune căutarea traducerii t (în limba țintă) care maximizează o combinație liniară a probabilităților diverșilor factori utilizați. Probabilitatea de traducere și regula de decizie este dată de:

$$t^* = \arg \max_e \sum_{k=1}^n \lambda_k h_k(t, s) \quad (1)$$

unde $h_k(t, s)$ este unul din cei n factori (o funcție de trăsături caracteristice perechii $\langle t, s \rangle$) iar λ_k este ponderea acestuia. Cei mai importanți factori care contribuie la probabilitatea de traducere provin din modelele de traducere, modelele de limbă, modelele de generare și cele de distorsiune. Combinarea modelelor de traducere din corpusuri comparabile și a modelelor de traducere din corpusuri paralele este dată de ponderile λ_k folosite pentru fiecare factor.

Stabilirea ponderilor λ_k se realizează prin proceduri de învățare automată. Una dintre cele mai utilizate proceduri este MERT (Minimal Error Rate Training) (Och, 2003) care presupune existența unor traduceri de referință, în raport cu care se induc valorile λ_k pentru care diferența dintre textul de tradus și textul de referință este minimă în raport cu o măsură de similaritate. Algoritmul MERT este un proces iterativ, computațional foarte intens, dar complet nesupervizat. Decodorul MOSES este însoțit în distribuția standard și de programul care implementează algoritmul MERT care implicit utilizează scorul BLEU (Papineni, et al., 2002). Această implementare (Bertoldi et al, 2009) a algoritmului MERT poate folosi opțional și alte măsuri de evaluare.

Pentru colectarea corpusurilor necesare proiectelor noastre a fost construit un program specializat care va fi descris în continuare. După cum se va vedea din prezentarea sistemului, acesta poate colecta atât corpusuri paralele cât și corpusuri comparabile, în funcție de natura conținutului multilingv al site-urilor prelucrate.

3. Descrierea sistemului de colectare a corpusurilor multilingve

O componentă importantă a motoarelor de căutare este un modul (numit *spider* sau *crawler*) care citește fiecare pagină și reține legăturile care pleacă din aceasta (Kobayashi and Takeda, 2000). Astfel se generează un graf plecând de la o listă inițială de adrese ce conțin referințe către alte locații. Procesul de generare poate fi automat (prin liste de căutare generate la prima indexare urmând ca la anumite intervale de timp să se verifice schimbări în conținutul acestora), manual (lista de site-uri este actualizată prin intervenție umană) sau hibrid.

Problema care apare în cazul extragerii de corpus este relevanța rezultatelor. În multe situații se obțin și referințe către site-uri fără nici o legătură cu tema de la care s-a plecat și pentru care nu se mai pot aplica aceleași reguli de indexare.

De exemplu, pe site-ul Wikitravel la articolul București se regăsesc următoarele referințe ce nu prezintă interes pentru extragerea de corpus, deoarece paginile respective nu se încadrează în tiparul acceptat de aplicație:

- **Carsrent** – există doar versiunea în limba engleză
- **Allrent, Budget** s – linkurile-ul sunt invalide

Din acest motiv, listele trebuie supuse unui proces de filtrare în funcție de specificul fiecărei pagini. Stabilirea link-urilor ce vor fi verificate se face prin filtrarea acestora conform unor criterii bine alese. De exemplu:

- Dacă se dorește ca articolul/documentul vizitat să facă parte din pagina curentă se poate recurge la eliminarea elementelor *href* ce conțin *http://* în interior (majoritatea referințelor interne se fac prin calea relativă către document) sau eventuala eliminare a legăturilor ce nu conțin adresa de bază a site-ului curent.
- Există cazuri în care se poate face o selectare bazată pe organizarea internă a paginii web. De exemplu, toate paginile de știri conțin în adresele lor */news/* iar toate paginile referitoare la vreme conțin în adresele lor */weather/*

Indiferent de metoda care se alege, dacă structura sitului se modifică la un moment dat, aplicația trebuie actualizată corespunzător, atrăgând după sine modificări în codul sursă, care devine greu de citit și reutilizat în alte aplicații.

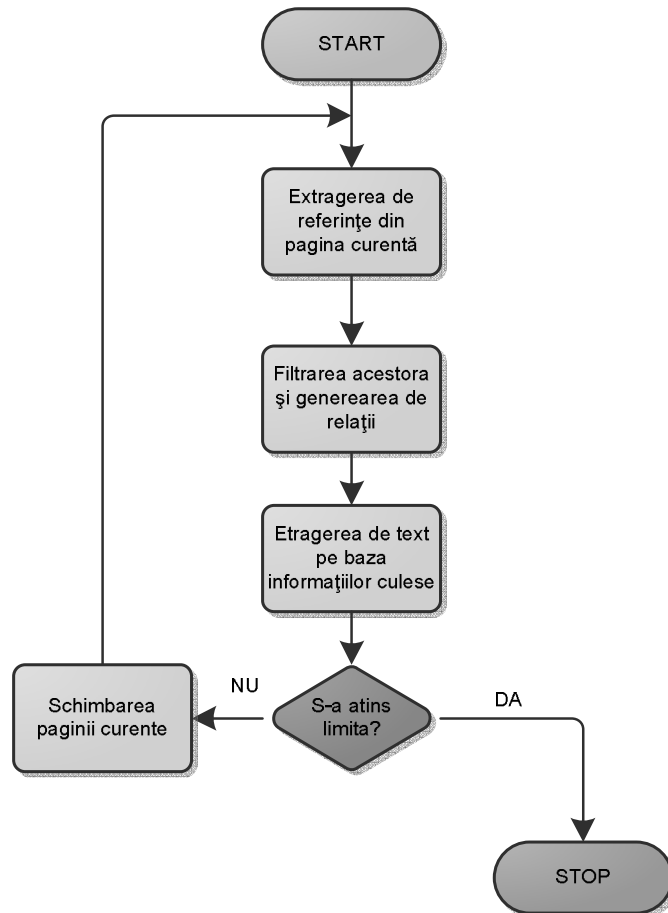


Figura 1: Diagrama de flux pentru algoritmul de extragere de text din pagini web

În continuare propunem o soluție de simplificare și minimizare a efortului depus la adăugarea unor noi criterii de indexare sau modificarea celor deja existente, prin utilizarea unei arhitecturi bazate pe diagrame de flux. În (Chappell, 2009) sunt prezentate avantajele unei soluții de acest tip. Astfel, fiecare bloc de procesare al diagramei de flux este privit independent de celelalte putând fi editat în orice moment și permițându-se adăugarea de pași suplimentari la fluxul deja existent. În fundal (eng. *background*) sunt lansate procese de tip consolă care primesc ca parametri de intrare fluxul și generează ieșire de tip text. Aceasta este preluată și este transmisă, în urma aplicării unei alte serii de prelucrări, la modulul următor.

În varianta actuală a sistemului am implementat două tipuri de blocuri, ilustrate în Figura 1, acestea fiind suficiente pentru operațiile uzuale într-un flux de date:

- blocuri de procesare – module ce vor genera sau altera textul
- blocuri decizionale – stabilesc drumul care va fi urmat

3.1 Arhitectura sistemului

Sistemul de colectare de corpusuri multilingve este compus din două module distincte: editorul de fluxuri și serviciul de colectare a paginilor web.

EDITORUL DE FLUXURI. În cadrul editorului de fluxuri se realizează schema de funcționare a aplicației, iar editarea diagramei poate fi făcută în mod vizual. De asemenea, pot fi editați parametrii lansării în execuție a proceselor ce vor alcătui sistemul de extragere, iar blocurile pot fi testate în mod individual.

SERVICIUL DE COLECTARE A PAGINILOR WEB este un serviciu Windows care rulează la un interval de timp prestabilit. Parametrii de configurare a serviciului sunt încărcăți dintr-un fișier XML. Documentul XML reprezintă schema de funcționare a aplicației creată cu ajutorul editorului de fluxuri. Serviciul de colectare a paginilor web, conform diagramei, lansează în execuție procesele și asigură comunicarea între ele. De asemenea, serviciul dispune de un jurnal detaliat al operațiilor efectuate cu ajutorul căruia pot fi investigate posibilele erori ale proceselor lansate.

3.2 Modul de utilizare

Pentru colectarea de corpus multilingv dintr-un anumit domeniu al web-ului este necesară crearea unei diagrame cu ajutorul aplicației de editare de fluxuri. În această diagramă sunt specificate procesele ce urmează a rula. De asemenea, în diagramă trebuie stabilite parametrii de lansare în execuție a proceselor (interpretorul utilizat, parametrii în linie de comandă etc.) și ce expresii regulate vor fi aplicate intrării/ieșirii proceselor. Pentru algoritmi în care este luată o decizie, în diagramă se va completa și regula de validare a conținutului.

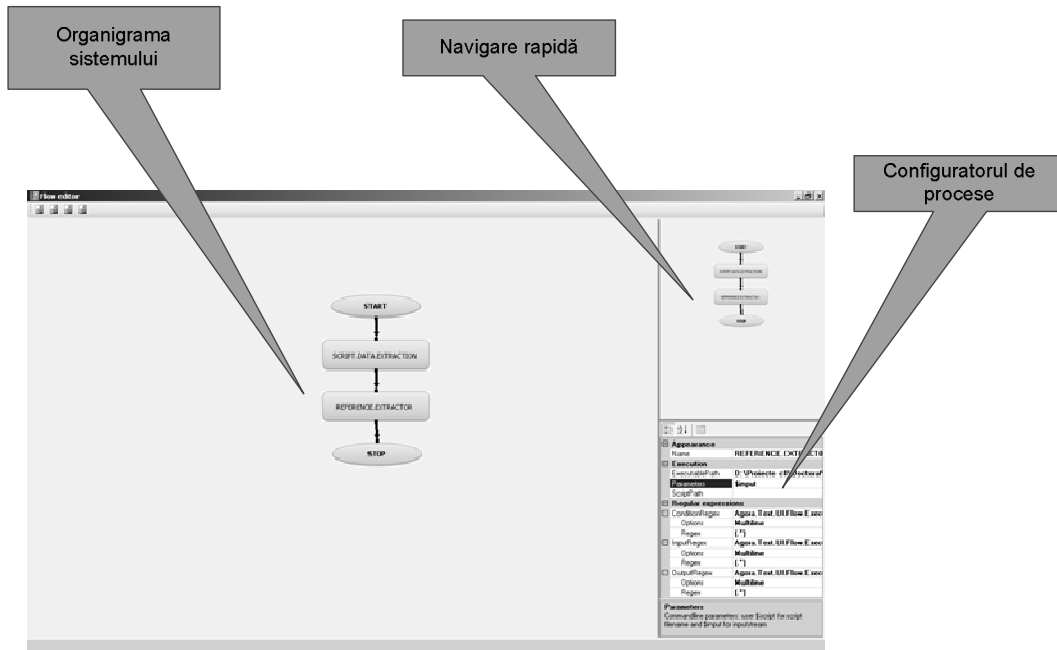


Figura 2: Editorul de diagrame

Fiecare bloc din diagramă reprezintă o unitate de execuție. Aceasta este caracterizată de calea către numele programului ce urmează a fi lansat, parametrii liniei de comandă (cuvintele cheie \$script și \$input vor fi înlocuite cu numele script-ului și, respectiv, intrarea obținută la pasul anterior) și expresiile regulate folosite (pentru intrare, ieșire și condiția de validare).

Odată definită, diagrama de flux este salvată într-un fișier XML ce este importat de sistemul propriu-zis de extragere. Acesta execută întregul flux la intervalul de timp stabilit în prealabil. Diagrama poate fi modificată în orice moment, urmând ca aceasta să fie reîncărcată la următoarea pornire a aplicației de colectare de pagini web.

4. Cazuri practice

Prezentăm în continuare două exemple de utilizare a aplicației: (i) pentru constituirea de corpus multilingv din domeniul turismului (*Wikitravel*) și pentru colectarea de corpus puternic comparabil din domeniul jurnalistic (*Parlamentul European* - secțiunea de știri). În ambele cazuri s-au folosit scheme de procesare simple similare cu cea prezentată în Figura 2. Au fost utilizate doar două programe/procese și nu au fost necesare blocuri decizionale. Cele două procese constau în: (i) generarea unei liste de adrese ce urmează a fi scanate și (ii) parcurgerea acestora și prelucrarea rezultatelor.

4.1 Wikitravel

În cazul site-ului Wikitravel s-a urmărit extragerea de corpusuri puternic comparabile din registrul textual al turismului pentru limbile română, engleză și germană. Cu ajutorul unei expresii regulate au fost extrase adresele paginilor din limba română

CONSTRUCȚIA AUTOMATĂ DE CORPUSURI COMPARABILE MULTILINGUALE

disponibile în domeniul „wikitravel.org”. Au fost urmate link-urile ce se regăseau în versiunea pentru limba română a paginii, urmând ca traducerile să fie găsite prin căutarea textului în subdomeniile /de/ sau /en/. Pentru validarea adreselor candidate au fost folosite liste pentru secvențele de adresă permise și pentru secvențele de adresă nepermise, prezentate în tabelul 1, prioritatea pentru validare fiind acordată secvențelor permise.

Tabel 1: Reguli de validare a adreselor

Secvențe de adresă permise	Secvențe de adresă nepermise
wikitravel.org/de/	#
wikitravel.org/en/	mailto
wikitravel.org/ro/	special
	action=
	Special:
	doubleclick
	User:
	Wikitravel:
	Utilizator:
	Image:
	http:
	www.

Legăturile care reprezentau traduceri ale articolului pentru engleză și germană au fost grupate, iar în cazul în care acestea lipseau, au fost înlocuite generic cu o pagină goală. Procesul a fost continuat prin extragerea textului și eliminarea elementelor HTML din date.

Tabel 2: Rezultatele obținute în cazul Wikitravel

română - engleză	română-germană
400.000 cuvinte	100.0 vint

4.2 Parlamentul European

Pe site-ul Parlamentului European la secțiunea știri se găsesc articole traduse în 22 de limbi. La fel ca în cazul Wikitravel, pornind de la adresa http://www.europarl.europa.eu/news/public/toute_actualite/default/default_ro.htm se parcurg toate referințele găsite pe pagină. Spre deosebire de cazul anterior, un site de știri se actualizează mai des, astfel că datele cantitative prezentate în Tabelul 3 sunt cele de la momentul scrierii acestui articol.

Tabel 3: Corpus obținut

Limbă	Cod	Cuvinte	Limbă	Cod	Cuvinte
bulgară	bg	54351	italiană	it	80431
cehă	cs	57480	lituaniană	lt	80431
daneză	da	67856	letonă	lv	67211

Limbă	Cod	Cuvinte
germană	de	70307
greacă	el	90729
engleză	en	84484
spaniolă	es	91447
estoniană	et	39388
finlandeză	fi	62692
franceză	fr	89171
ungară	hu	59584

Limbă	Cod	Cuvinte
malteză	mt	46041
olandeză	nl	74390
poloneză	pl	56031
portugheză	pt	75557
română	ro	81433
slovacă	sk	67189
slovenă	sv	68794

Sperăm ca în câteva luni, corpusul multilingv EU-News (22 limbi) să devină o resursă publică, de dimensiune corespunzătoare, pentru cercetări și dezvoltări în domeniul cros și multilingv (traducere automată, sisteme de regăsire documentară multilingvă, sisteme de întrebare/răspuns cros-lingvistice, etc.).

În ambele exemple prezentate, corpusurile au fost stocate în directoare distincte de fișiere de tip text. Documentele paralele/comparabile au fost numite folosind identificatori unici extensiile acestora denotând codul ISO al limbii documentului (de pildă nume de tipul xxxxxxxx.bg, xxxxxxxx.cs, ..., xxxxxxxx.sv identifică fișierele conținând documente paralele/comparabile în cele 22 de limbi ale corpusului multilingv). În plus, au fost generate metadate corespunzătoare fiecărui fișier text, metadate ce documentează calea către fișier, adresa de la care a fost extras fișierul, numărul de cuvinte din fișier etc.

5. Concluzii

Rezultatele bune obținute în urma colectării de corpus multilingv din site-ul Wikitravel și din secțiunea de știri a site-ului Parlamentului European dovedesc utilitatea setului de instrumente propus în această lucrare. Aceste instrumente facilitează interoperabilitatea diverselor aplicații scrise în diferite limbafe de programare, fiecare aplicație fiind reprezentată ca un proces al unei diagrame de flux. De asemenea, posibilitatea de a reconfigura în orice moment a diagrama de flux permite abordarea unor proiecte de colectare de corpus de complexitate ridicată – procesele, și regulile de validare pot fi schimbate chiar în cursul colectării documentelor.

Mulțumiri. Activitatea de cercetare descrisă a fost sprijinită de Comisia Europeană prin proiectul ACCURAT (248347/FP7) și (parțial) de CNCSIS – UEFISCSU prin proiectul PNII – IDEI STAR, 742/19.01.2009

Referințe bibliografice

- Bertoldi, N., Haddow, B., Fouet, J.-B. (2009). *Improved Minimum Error Rate Training in Moses*, în Prague Bulletin of Mathematical Linguistics, Nr. 91, 2009, pp. 7-16.
- Ceaușu, A. (2009). *Tehnici de traducere automată și aplicabilitatea lor limbii române ca limbă sursă*, București, România: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.

- Chappell, D. (2009). *The Workflow Way - Understanding Windows Workflow Foundation*.
- Germann, U. (ed.) (2001). *Aligned Hansards of the 36th Parliament of Canada - Release 2001-1a*. <http://www.isi.edu/natural-language/download/hansard/> (15.03.2010).
- Hoang, H., Koehn, P., Lopez, A. (2009). A Uniform Framework for Phrase-Based, Hierarchical and Syntax-Based Machine Translation, International Workshop on Machine Translation (IWSLT), pp. 152-159.
- Koehn, Ph. (2005). *EuroParl: A Parallel Corpus for Statistical Machine Translation*. Machine Translation Summit 2005. Phuket, Thailand. <http://www.statmt.org/europarl/> (16.03.2010).
- Koehn, Ph., Hoang, H., Birch, A., Callison-Burch, Ch., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, Ch., Zens, R., Dyer, Ch., Bojar, O., Constantin, A., Herbst, E. (2007). *Moses: Open Source Toolkit for Statistical Machine Translation*, Annual Meeting of the Association for Computational Linguistics (ACL), demonstration session, Prague, Czech Republic, June 2007.
- Kobayashi, M. and Takeda, K. (2000). „Information retrieval on the web”. *ACM Computing Surveys* (ACM Press) **32** (2): 144–173. doi:10.1145/358923.358934
- Och, F. J. (2003). *Minimal Error Rate Training in Statistical Machine Translation*. In Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics, July 2003, pp. 160-167.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W. J. (2002). *BLEU: a method for automatic evaluation of machine translation*. In ACL-2002: 40th Annual meeting of the Association for Computational Linguistics pp. 311–318.
- Steinberger, R., Pouliquen, B., Widiger, A., Ignat, C., Erjavec, T., Tufiș, D., Varga, D. (2006). *The JRC-Acquis: A Multilingual Aligned Parallel Corpus with 20+ Languages*. In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)*, pp. 2142-2147, Genoa, Italy, May 2006. ELRA - European Language Resources Association.
- Tiedemann, J. (2009). *News from OPUS - A Collection of Multilingual Parallel Corpora with Tools and Interfaces*. In N. Nicolov and K. Bontcheva and G. Angelova and R. Mitkov (eds.) *Recent Advances in Natural Language Processing* (vol V), pages 237-248, John Benjamins, Amsterdam/Philadelphia.
- Tufiș, D., Ion, R., Ceaușu, A., Ștefănescu, D. (2007). *Servicii web lingvistice ale ICIA*. În Ionuț Pistol, Dan Cristea, and Tufiș, D. (eds.), *Resurse Lingvistice și Instrumente pentru Prelucrarea Limbii Române*, pp. 61-68, Iași, România, dec. 2007. Universitatea "Al.I. Cuza" Iași, Editura Universității "Al.I. Cuza" Iași. ISBN 978-97-3703-297-3.
- Tufiș, D., Koeva, S., Erjavec, T., Gavrilidou, M., Krstev, C. (2009). Building Language Resources and Translation Models for Machine Translation Focused on South Slavic and Balkan Languages. In J. Machacova and K. Rohsmann (eds) „Scientific Results of the SEE-ERA.NET Pilot Joint Call”. Center for Social Innovation Publisher, Vienna, ISBN 978-3-200-01567-8, pp. 37-48, June 2009.

PARSAREA COMPARATIVĂ A DICȚIONARELOR-TEZAURE ROMÂNEȘTI, FRANCEZE ȘI GERMANE

NECULAI CURTEANU¹, MIHAI ALEX MORUZ^{1,2}, DIANA TRANDABĂȚ¹

¹*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

²*Universitatea „Al. I. Cuza”, Facultatea de Informatică, Iași – România;*

{curteanu, mmoruz, dtrandabat}@iit.tuiasi.ro

Rezumat

Lucrarea de față prezintă un studiu comparativ privind modelarea lexico-semantică și parsarea dicționarilor tezaure pentru diferite limbi. Transformarea eficientă a unui tezaur într-o bază de date structurată de arbori lexico-semantici se bazează pe metoda configurațiilor SCD (Segmentare-Coeziune-Dependență), acestea fiind aplicate succesiv pentru identificarea segmentelor lexicografice și extragerea arborelui de sensuri și sub-sensuri din intrare. Modelarea lexico-semantică dezvoltată pentru tipologia intrărilor de dicționar din Dicționarul tezaur al Limbii Române - format *nou* (**DLR**) a fost adaptată și aplicată la parsarea altor patru tezaure: Dicționarul tezaur al Limbii Române - formatul *vechi* (**DAR**), Trésor de la Langue Française (**TLF**), Deutsches Wörterbuch – GRIMM (**DWB**), și Göthe-Wörterbuch (**GWB**). Lucrarea prezintă clasele de markeri de sensuri și hipergraful de dependențe specifice pentru fiecare dicționar, parsarea acestor tezaure, și discuții privind rezultatele obținute și analiza erorilor.

1. Introducere

Parsarea unui dicționar presupune transformarea intrărilor ce conțin text sub formă de glosă, într-un format indexabil. Astfel, fiecare intrare de dicționar este transpusă într-o structură complexă, care conține atât sensurile definite, precum și descrieri detaliate ale formei intrării, cu referire la ortografie, morfologie, fonetică, etimologie, uz etc.

Scopul acestei lucrări este analiza impactului configurațiilor SCD (Segmentation-Cohesion-Dependency) (Curteanu; 2006) la parsarea intrărilor de dicționar pentru patru dicționare tezaure: Dicționarul tezaur al Limbii Române - formatul *vechi* (**DAR**), Trésor de la Langue Française (**TLF**), Deutsches Wörterbuch – GRIMM (**DWB**), și Göthe-Wörterbuch (**GWB**). Metoda a fost deja aplicată cu succes pentru parsarea Dicționarului tezaur al Limbii Române - formatul *nou* (**DLR**) (Curteanu et al.; 2008a, 2008b, 2009a), folosind configurații SCD. Prima configurație are rolul de a identifica segmentele lexicografice distincte dintr-o intrare de dicționar. În cadrul strategiei dezvoltate pentru **DLR**, alte două configurații diferite de parsare sunt folosite: una care identifică și extrage, pentru fiecare intrare din dicționar, arborele specific de sensuri (ierarhia sensurilor principale și secundare), și o altă configurație care parsează fiecare nod din arborele de sensuri cu scopul de a clasifica definițiile aceluși sens/subsens din dicționar. Spre deosebire de metodele standard de parsare a textului de dicționar, în care toate câmpurile unei intrări de dicționar sunt analizate secvențial, noua metodă reușește detașarea procesului de construire a arborelui de sensuri (a *doua configurație* SCD) de procesul parsării la definițiile sensurilor (a *treia configurație* SCD). Separarea celor

două procese se face în principal prin selectarea *breadth-first* a tuturor marcherilor la sensuri, urmată de analiza *depth-first* a secvențelor de marcheri pentru fiecare intrare de dicționar. Amintim că o *primă configurație* SCD realizează recunoașterea segmentelor lexicografice ale intrării.

Spre deosebire de orientările standard în parsarea intrărilor de dicționar (Neff & Boguraev; 1989), de exemplu sistemul **LexParse** (Hauser & Storrer; 1993), (Kammerer; 2000), (Lemnitzer & Kunze; 2005) sau gramaticile lexicografice (Curteanu & Amihăesei; 2004), (Tufiș; 2001), metoda descrisă în (Curteanu et al., 2008a) *detașează complet* procesul de construire a arborilor de sensuri de procesul parsării definițiilor la sensuri. O configurație SCD are următoarele componente:

- Un set de clase de marcheri: un marcher reprezintă o graniță computațional-lingvistică pentru o categorie lingvistică specifică;
- O ierarhie de tip arbore (sau graf aciclic orientat, sau hipergraf), care stabilește dependențele dintre clasele de marcheri;
- Un algoritm de parsare, care execută următorii pași: recunoașterea marcherilor, identificarea structurilor textuale dintre doi sau mai mulți marcheri, și clasificarea acestor structuri ținând cont de ierarhia claselor de marcheri. Pentru aplicarea configurațiilor SCD la parsarea dicționarelor este necesară o abordare lexico-semantică a claselor de marcheri de sensuri și definiții, și stabilirea unei reprezentări ierarhice a lor.

Lucrarea de față prezintă configurațiile SCD pentru cele patru dicționare-tezaur considerate:

- *Recunoașterea segmentelor lexicografice care corespund unei intrări de dicționar:* Pentru dicționare complexe / tezaure, fiecare intrare este compusă din mai multe segmente lexicografice, întinderi textuale care partiționează intrarea de dicționar în unități lexicografice distincte. Astfel, o *primă configurație* SCD este dedicată obținerii segmentelor lexicografice specifice unei intrări.

- *Construcția arborelui de sensuri al unei intrări:* După recunoașterea segmentelor lexicografice, acestea sunt clasificate în funcție de structura lor lingvistică și semantică. Cel mai important segment este cel al *descrierii sensurilor*, asupra căruia se aplică o *a doua configurație* SCD pentru extragerea arborelui de sensuri al intrării. Pentru **DLR**, arborele de sensuri este obținut cu o precizie de 91.18% (Curteanu et al.; 2008a, 2009a).

- *Clasele de marcheri* folosite pentru identificarea arborelui de sensuri al unei intrări din **DLR** sunt, pentru *sensuri principale*: majuscule (**A.**, **B.**, etc.); cifre romane (**I.**, **II.**, etc.); cifre arabe (**1.**, **2.**, etc.). Pentru *sensuri secundare*: rombul plin **◆**; rombul gol **◇**; enumerarea literală cu litere latine mici **a)**, **b)**, **c)**,...

- *Segmentarea definițiilor:* A *treia configurație* SCD folosește un set specific de clase de marcheri pentru segmentarea la definiții atomice în fiecare sens al intrării și stabilirea dependențelor dintre sensuri (vezi (Curteanu et al.; 2009a) pentru o clasificare a definițiilor atomice). Analiza intrărilor **DLR** a dus la următoarele tipuri de definiții: **(i)** *MorfDefs* – definiții morfologice; **(ii)** *RegDefs* – definiții scrise cu font regular; **(iii)** *BoldDefs* – definiții scrise cu bold; **(iv)** *ItalDefs* – definiții scrise cu italic; **(v)** *SpecDefs* – definiții ce conțin specificații; **(vi)** *SpSpecDefs* – definiții scrise cu litere spațiate, ce conțin anumite specificații; **(vii)** *DefExems* – exemple la definiții, cu rolul de a întregi înțelesurile unei definiții. Tipurile de definiții propuse aici primesc roluri funcționale specifice în

descrierea sensurilor principale, secundare, sau de granularitate semantică mai fină (Curteanu et al.; 2009a).

2. Parsarea Dicționarului tezaur al limbii române– formatul vechi (DAR)

Structura principalelor segmente lexicografice ale unei intrări **DAR** este următoarea:

I. Segmentul *traducerii franceze*, *FreSeg*, conține lema și structura principalelor sensuri ale cuvântului, traduse în limba franceză. Comparația dintre traducerea structurii sensurilor lemei și descrierea completă a intrării **DAR** pentru limba română revine la aplicarea operației de *subsumare* între arborele de sensuri al cuvântului românesc și arborele de sensuri al traducerii lui în limba franceză. Există multiple situații în această comparare, inclusiv aceea în care perechea <*cuvânt*_{Rom}, *mot*_{Fre}> nu poate fi determinată: de exemplu <*mămăligă*_{Rom}, ?_{Fre}>.

II. Segmentul *descrierii generale* a lemei din limba română (*RomSeg*). *RomSeg* conține, uneori pe mai multe *paragrafe*, informații morfo-sintactice, semantice, de utilizare etc. ale intrării. *RomSeg* începe, de obicei, cu lema în format italic (uneori lema se află pe primul rând al primului paragraf din *RomSeg*).

III. Al treilea segment al unei intrări din **DAR**, *SenseSeg*, este chiar *descrierea lexical-semantică* a cuvântului-lemă. *SenseSeg* constituie principalul obiectiv al analizei lexicosemantice și al programului de parsare la arborele de sensuri în **DAR**.

IV. Al patrulea segment, *NestSeg*, al unei intrări în **DAR** este format din unul sau mai multe „*cuiburi*”, întinderi de text lexical-semantic în care sunt descrise variante morfologice-sintactice, fonologice, regionalisme etc. ale intrării, uneori cu descrierea subsidiară a sensurilor specifice. Structura segmentului-*cuib* din **DAR** este similară cu cea a unei intrări generale din **DAR**. Construcția posibil recursivă a intrărilor **DAR** rezultă din necesitatea parsării la sensuri și într-un segment-*cuib*.

V. Al cincilea segment al intrărilor din **DAR**, denotat *EtymSeg*, face precizări asupra etimologiei lemei și este introdus printr-o *liniuță-de-etimologie* (cratimă-mare „-”).

Din cele cinci segmente ale unei intrări **DAR**, singurele obligatorii sunt *FreSeg* și *SenseSeg*. Celelalte trei segmente apar opțional în descrierea intrărilor, în funcție de necesitățile și specificul fiecărei leme.

2.1 Clase de marcheri la sensuri

În ordinea priorităților din ierarhia claselor de marcheri la sensuri **DAR**, avem:

1. Litere latine majuscule (*LatCapLett_Enum*): A., B., C., ... 2. Cifre romane majuscule (*LatCapNumb_Enum*): I., II., ... 3. Cifre arabe (*ArabNumb_Enum*): 1⁰., 2⁰., ...

În mod similar cu clasificarea sensurilor din **DLR**, vom considera că marcherii acestor clase introduc *sensuri principale* și în **DAR**.

4. Pentru a introduce *sensuri secundare*, de obicei pentru definițiile *obligatorii* și/sau *opționale* (Curteanu et al.; 2009a), **DAR** folosește marcheri de sens care introduc și în **DLR** definițiile de tip *MorfDef*, *RegDef*, *BoldDef*, *ItalDef*, *SpecDef*, *SpSpecDef*, și *DefExem*, dar și câțiva *marcheri specifici DAR*: ||, |, #, †. **DAR** definește primii doi

marcheri astfel: „Liniuțele verticale despart diferitele subsensuri, și anume liniuța simplă | pe cele mai puțin deosebite, iar liniuța dublă || pe cele mai marcante. Uneori, la verbe, aceste liniuțe despart pe *transitive*, de *intransitive* și de *reflexive*”. Marcherul # delimitează „idiotismele și locuțiunile proverbiale”, iar † introduce sensurile sau expresiile scoase din uzul limbii (Curteanu et al.; 2009b).

5. În funcție de nivelul din arborele de sensuri al descrierii lexical-semantice, situându-se fie pe sensuri principale fie secundare, **DAR** folosește și *enumerarea literală* pe două niveluri: (5.a) Litere latine minuscule (*LatSmaLett_Enum*): a.), b.), ... (5.b) Sub o enumerare literală *LatSmaLett_Enum* se poate face o a doua enumerare, cu litere grecești minuscule (*GreSmaLett_Enum*): α.), β.), γ.), ...

2.2 Ierarhia claselor de marcheri pentru dependențele la sensuri în DAR

Introducem următoarele ierarhii și subierarhii pentru clasele de marcheri la sensuri și definiții din **DAR**:

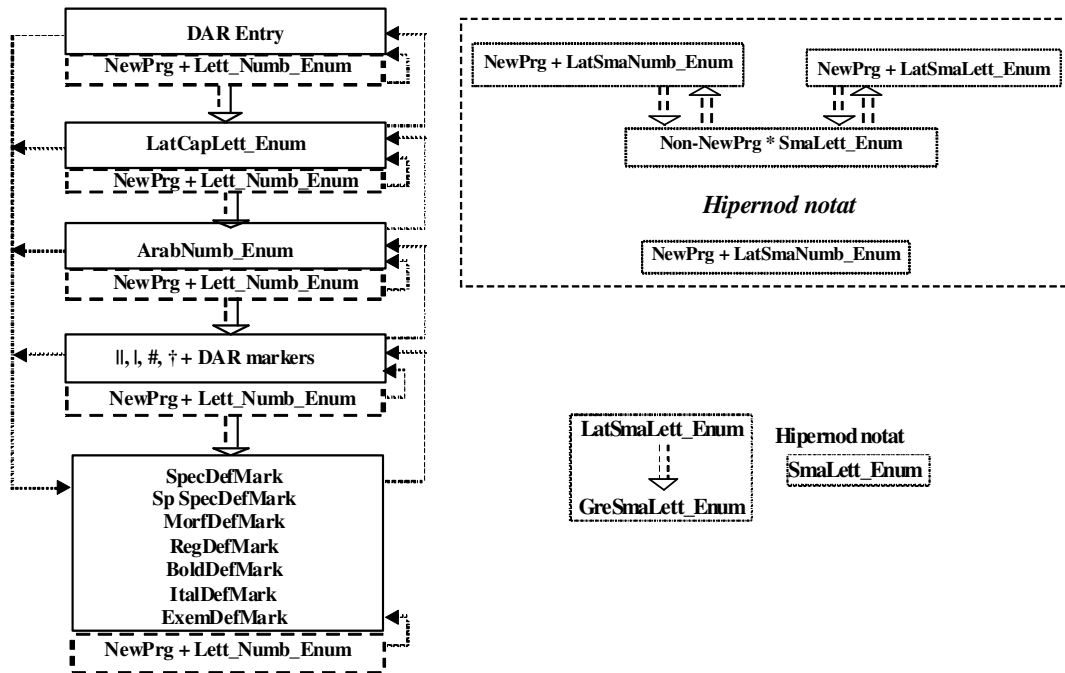


Figura 1: Hipergraful de dependențe între clasele de marcheri la sensuri în DAR

Hipergraful de dependențe la clase de marcheri **DAR** (Curteanu et. al.; 2009b) este mult mai complex (enumerări la paragrafe noi și enumerări literale) în comparație cu hipergraful pentru **DLR** (Curteanu; 2008a, 2009a).

2.3 Probleme speciale și rezultate ale parsării DAR

Principala *dificultate* a parsării intrărilor din **DAR** este folosirea marcherului de sens *Paragraf_Nou* (*NewPrg*), aspect care *nu* apare în **DLR**. Marcherul *NewPrg* este costisitor de recunoscut (din punct de vedere al timpului de execuție) deoarece, pentru acest marcher, a trebuit să introducem un *nivel implicit* de *enumerare* dacă și numai dacă *NewPrg nu* este însoțit și de un alt tip de marcher la sensuri **DAR** (în special

marcheri de sensuri principale). Enumerarea implicită generată de marcherul *NewPrg* se face cu *cifre latine minuscule* (*LatSmaNumb_Enum*): *i, ii, iii, iv, v, vi, vii* etc.

A doua problemă dificilă care apare la parsarea **DAR** este procedeul de rafinare a marcherului *NewPrg* prin enumerare literală cu latine minuscule *LatSmaLett_Enum*, care pot, la rândul lor, să fie rafinate printr-o enumerare implicită de tipul *NewPrg*. Această problemă a fost rezolvată atribuind marcherilor de tip *NewPrg*, și enumerării cu latine minuscule, niveluri de sens care variază în funcție de context.

Un exemplu remarcabil este intrarea **CAL** din **DAR**, în care aspectul ce trebuie subliniat este că *NewPrg*, enumerat cu cifre latine minuscule *LatSmaNumb_Enum* (pe primul subnivel de dependență), apelează enumerarea literală cu *NewPrg* (la paragraf nou) – al doilea subnivel de dependență, care la rândul ei apelează *NewPrg* enumerat (al treilea subnivel). De pe acest nivel de *NewPrg* este apelată din nou enumerarea literală la *Non-NewPrg* (al patrulea subnivel de dependență). Toate cele *cinci* niveluri componente de dependență se atașează formal nivelului de sens reprezentat de marcherul 1⁰. Aceste apelări reciproce și succesive se bazează pe hipergraful de dependențe din Fig. 1.

Parsarea **DAR** este descrisă în detaliu în (Curteanu et al.; 2009b). Au fost parsate un număr de 37 intrări **DAR**, de dimensiuni mari și medii. Evaluarea manuală a parsării este foarte promițătoare având în vedere problemele ridicate de parsarea **DAR**, semnificativ mai dificile decât la **DLR**. În lipsa unui corpus gold pentru parsarea intrărilor **DAR**, nu s-a putut face o evaluare automată a parserului.

3. Parsarea TLF cu configurații SCD

3.1 Segmente lexicografice în TLF

Dicționarul-tezaur **TLF** prezintă și asemănări cu **DLR** dar și aspecte distincte interesante. Ca observație generală, **TLF** este unul din cel mai bine structurate dicționare-tezaur dintre cele analizate. Configurația *SCD-Config1* realizează diviziunea la *segmente lexicografice* a intrărilor din **TLF**, structura acestor segmente fiind relativ simplă. O intrare **TLF** începe chiar cu segmentul de *descriere a sensurilor*. Acest segment este urmat, opțional, de o serie de segmente *finale* introduse de etichete **TLF** bine definite. Aceste segmente pot fi văzute și ca un singur segment lexicografic, format din componente opționale, specifice fiecărei etichete. Calupul de segmente finale (sau segmentul final) este introdus de etichete după modelul ce urmează:

REM. 1. ... 2. ... 3. ...
PRONONC. ET ORTH. – ... Eng.: ...
ÉTYMOL. ET HIST. I. ... 1. a) ... b) ... 2. ... 3. ... II. ...
STAT. Fréq. abs. littér.: ... Fréq. rel. littér.: ...
DÉR. 1. ... 2. ... 3. a) ... b) Rem. a) ... b) ...
BBG. ...

Există intrări **TLF** în care aceste segmente componente nu încep la paragraf nou *NewPrg* ci formează împreună un paragraf compact, segmentele componente fiind delimitate de etichetele de mai sus (în unele cazuri cu *majuscule-bold*, în alte cazuri cu *minuscule-bold*, sau amestecat). După cum se observă din șablonul de mai sus, unele

segmente finale conțin descrieri particulare de senzori, perfect similare cu descrierile din segmentul propriu-zis de senzori.

3.2 Clase de marcheri TLF și hipergraful de dependență

Marcherii de senzori din TLF sunt oarecum asemănători cu cei din DLR dar există și diferențe semnificative, ce impun până la urmă specificul de construcție și de reprezentare al celor două dicționare-tezaur. În comparație cu DLR, în TLF avem:

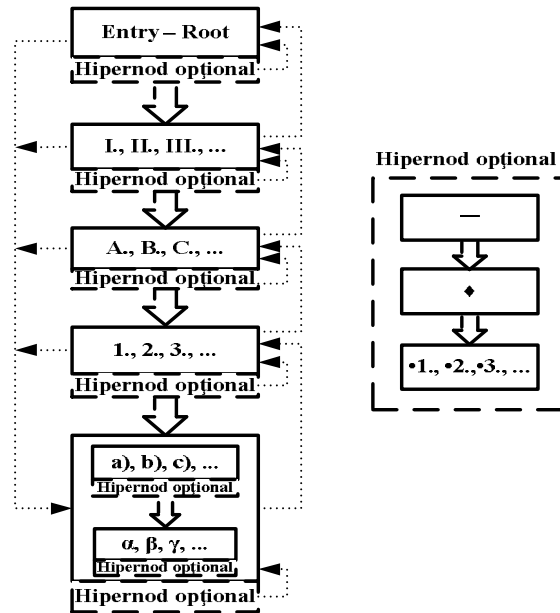


Figura 2: Hipergrafurile de dependențe la clasele de marcheri TLF

(1) *Marcheri noi:* (1.1) Marcherul „-” (*liniuța-de-moștenire*) introduce un *subsens moștenit*. Pentru marcherul „-” de moștenire a unui sens avem *două* situații importante în intrările TLF: (a) Atunci când „-” apare *după* un alt marcher de sens din TLF, rolul „-” este cel de moștenire a sensului de pe nodul-părinte de sens (*regent* sau *nu*) din arborele de senzori al intrării; (b) Când „-” este la *paragraf nou*, deci *NewPrg* + „-”, liniuța-de-moștenire „-” are rol de *subsens intermediar*, similar cu rolul *NewPrg* din tezaurul DAR. În acest caz particular, „-” introduce un subsens al marcherului / sensului-părinte *regent*, *direct* dependent de (sub)sensul regent. (1.2) Marcherul •1., •2., ... (*bulină-roșie indexată*) definește un *concept nou, specific TLF: DefExems indexate* pe întreaga intrare TLF (*IdxDefExem*). (1.3) *Enumerarea literală dublă:* enumerarea *LatSmaLett_Enum* (latine mici a), b), ...) este rafinată prin *GreSmaLett_Enum* (litere grecești mici α), β), ...).

(2) În TLF întâlnim ca marcher de sens secundar numai *rombul-plin* „♦”, nu și *rombul-gol*, prezent în DLR.

(3) În comparație cu DLR, în TLF apare o *inversare* a dependențelor la clasele de marcheri pentru senzurile principale, reflectată în hipergraful de dependențe din Fig. 2.

3.3 Rezultate de parsare la sensuri pentru marcheri expliciți în TLF

Pentru TLF au fost parsate un număr de 31 intrări, cu dimensiuni medii și mari (în medie, 5 pagini format A4). Rezultatele parsării TLF au fost evaluate manual, estimarea preciziei parsării fiind foarte bună. Principalul motiv al acestor rezultate este elaborarea consecventă a intrărilor TLF, structura lor regulată, apariția foarte rară a excepțiilor de la ierarhia de marcheri propusă. În lipsa unui corpus gold pentru parsarea intrărilor TLF, nu s-a putut face o evaluare automată (precisă) a parserului.

4. Segmente lexicografice și clase de marcheri la sensuri în DWB și GWB

4.1 Moduri de recunoaștere a segmentelor lexicografice din DWB

Intrările DWB sunt constituite din mai multe *segmente lexicografice*. Aceste segmente au o structură neuniformă și neunitară. *Caracteristic* pentru tezaurele germane DWB și GWB este structura segmentelor lexicografice, formate dintr-un prim subsegment (opțional) pe care l-am numit *sensul-rădăcină* al segmentului, și din al doilea subsegment, *corpul propriu-zis* al segmentului lexicografic. În timp ce sensurile incluse în *corpul* fiecărui segment lexicografic sunt descrise cu marcheri expliciți, relativ ușor de recunoscut, parsarea completă a segmentelor DWB nu este o sarcină ușoară. În DWB, segmentele lexicografice sunt introduse în *trei* moduri:

(A) După *sensul-rădăcină* (*definiția-rădăcină*) al unei intrări DWB, sau după *sensul-rădăcină* al unui segment lexicografic, sunt poziționate un număr de cuvinte-rezervate, scrise (de obicei) *spatiat* și cu *font-italic*. Aceste cuvinte-rezervate (sau *listă* de cuvinte rezervate) se constituie în *eticheta segmentului lexicografic* a cărui întindere **urmează** segmentului curent din DWB. Exemple de asemenea cuvinte-cheie ce delimitează începutul unui segment lexicografic: „*Form, ausbildung und ursprung*”, „*Formen*”, „*Ableitungen*”, „*Verwandtschaft*”, „*Verwandtschaft und form*”, „*Formelles und etymologisches*”, „*Gebrauch*”, „*Herkunft*”, „*Grammatisches*”, etc., sau, pentru segmentul de descriere a sensurilor din DWB, „*Bedeutung und gebrauch*” (sau numai „*Bedeutung*”). În exemple, etichetele de segmente lexicografice au fondul 25% gri:

GRUND, *m., dialektisch auch f. gemeingerm. wort; fraglich ist das geschlecht von got. *grundus in grunduwaddjus, vgl. afgrundiþa; sonst meist masc.: ahd. grunt, crunt; mhd. grunt; as. meist f., selten m.; grunte f., lett. grunts m., grunte f., poln. russ. slov. nlaus. grunt m. form und herkunft.*

1) für das verständnis der vorgeschichte des wortes ist die *z w i e g e s c h l e c h t i g k e i t* ...

H. V. SACHSENHEIM *spiegel* 177, 30;

die neuen grundt zu der kirchen *zimm. chron.*² 2, 539, 36; du findest noch vil gar alter meür und grunt und thürn SIGMUND MEISTERLIN in *städtechron.* 3, 51, 14. drey starcke grund 6, 290. *b e d e u t u n g . die bedeutungsgeschichte des wortes lässt sich schwer aufbauen, doch erschöpft diese unterscheidung einer mehr räumlichen und mehr flächenhaften vorstellung die sache nicht.*

I. grund bezeichnet die feste untere begrenzung eines dinges.

A. grund von gewässern; seit ältester zeit belegbar: profundum (sc. mare) crunt *ahd. gl.* 1, 232, 18; latid thea odra (*fisch*) eft an gr. faran *Hel.* 2633.

1) am häufigsten vom meer (in übereinstimmung mit dem anord. gebrauch): ...

(B) Al doilea mod de specificare a segmentelor lexicografice din DWB este cel prin care *imediat după* marcherii de sensuri principale sunt specificate cuvintele-cheie ce reprezintă *eticheta* segmentului lexicografic care urmează.

GEBEN, *dare*.

I. Formen, ableitungen, verwandschaft.

1) *es ist ein allgemein, aber ausschliesslich germanisches wort: goth. giban (praet. gaf), ahd. ...*

II. Bedeutung und gebrauch.

1) *geben und nehmen, die beiden sich ergänzenden gegenstücke, verdienen die erste...*

(C) Al treilea mod de specificare a descrierii lexical-semantice a unei intrări în **DWB** (deci și a segmentului de descriere a sensurilor), care este și cel mai frecvent, este acela în care nu este folosită nicio etichetă pentru începutul marcherilor de sensuri. Implicit, după sensul rădăcină (care poate fi redus doar la traducerea latină a cuvântului leamă din **DWB**), urmează segmentul de descriere a sensurilor (care este și unic), *fără eticheta explicită „Bedeutung”*, doar prin prezența expresă a marcherilor de sensuri, a definițiilor atomice, a exemplelor la definiții și a siglelor pentru aceste exemple.

4.2 Clasele de marcheri și dependențele la sensuri în **DWB** și **GWB**

Marcherii la sensuri sunt cei obișnuiți din dicționarele-tezaur, cu mențiunea că rafinarea sensurilor prin *enumerare literală* se face în **DWB** pe trei niveluri, *LatSmaLett_Enum*, *GreSmaLett_Enum*, dar și *Gre2SmaLett_Enum* ($\alpha\alpha$), ($\beta\beta$), (...), iar în **GWB** avem doar primele două. **GWB** conține și marcherul *liniuța-de-moștenire* (ca în **TLF**). Hipergrafurile de dependențe la sensuri în **DWB** și **GWB** sunt în Fig.3. și Fig.4.:

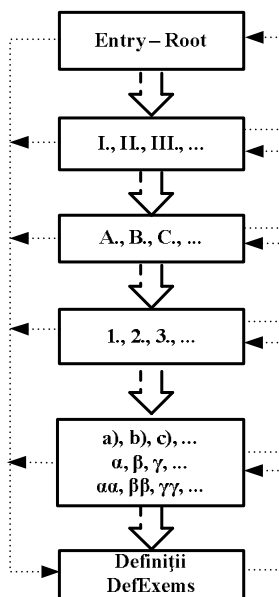


Figura 3: Hipergraful de dependențe a claselor de marcheri **DWB**

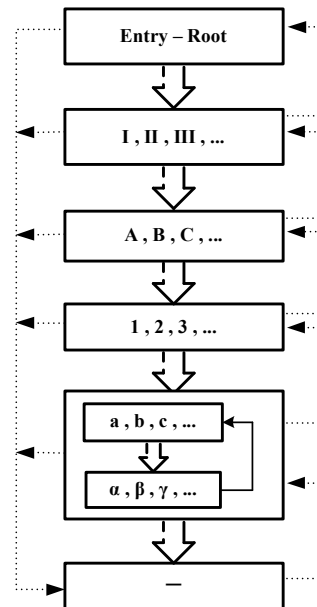


Figura 4: Hipergraful de dependențe a claselor de marcheri **GWB**

4.3 Rezultate de parsare la segmente lexicografice în **DWB**

Au fost parsate un număr de 17 intrări de mari dimensiuni din **DWB**, cu diverse particularități de redactare. La evaluarea manuală ne-am concentrat asupra rezultatelor de parsare a primelor două configurații SCD. Parsarea segmentului de descriere de sensuri se face cu mare precizie, însă separarea segmentelor lexicografice este mai puțin precisă, în mare măsură din cauza structurii complexe a segmentelor lexicografice în

DWB și a variației foarte mari a modului în care acestea sunt marcate. În lipsa unui corpus gold, nu s-a putut face o evaluare automată precisă a parserului.

5. Discuții

Este important să punem în evidență, comparativ lingvistic, rolul funcțional al *structurilor* lexico-semantice, *de sens*, cu cel al *marcherilor* la aceste structuri în dicționarele-tezaure analizate. Noile denumiri propuse au scopul de a stabili o mai strânsă legătură între sensuri și marcherii tipografici / lexicografici prin care sunt identificate și relaționate aceste sensuri.

MorfDef ⇒ *MorfDef*; *RegDef* ⇒ *GlossDef*; *BoldDef* ⇒ *IdiomDef*; *ItalDef* ⇒ *CollocDef*; *DefExem* ⇒ *DefExem*; *SpecDef* ⇒ *SpecDef*; *SpSpecDef* ⇒ *SpSpecDef*; *IdxDefExem* (concept nou: *DefExem* indexat).

Câteva argumente care susțin propunerile de mai sus: (1) Perechea *RegDefs–DefExems* este reprezentată în **DLR** (ca și în **DAR** și **TLF**) prin perechea de marcheri-fonturi *Regular–Italic*, în timp ce aceeași pereche de noțiuni este denotată în tezaurele **DWB** și **GWB** prin perechea *inversă* de marcheri-fonturi: *Italic–Regular*. Deci concepte similare sau identice pot fi codificate, în tezaure diferite, prin marcheri tipografici distincți. (2) Anumite noțiuni lexico-semantice pot fi particulare unui anumit dicționar-tezaur; așa cum este *DefExem Indexat* (codificat *IdxDefExem*), introdus în **TLF** printr-o bulină roșie indexată. *IdxDefExem* este atașat unei anumite definiții de sens, dar acest *DefExem* particular poate fi referit în mod specific pentru intrarea **TLF** în care este definit.

Experiența acestei lucrări a demonstrat încă o dată generalitatea și eficiența strategiei de parsare cu *configurații* SCD, în acest caz, a textelor de dicționare-tezaure complexe, fiecare cu multe caracteristici specifice. Principala realizare, tradusă în randamentul foarte bun al programelor de parsare, este aceea că analiza se face în primul rând pe secvențele de marcheri extrași din textul intrării, iar procesele de obținere a arborelui de sensuri și de parsare la sensuri / definiții atomice pot fi complet separate (Curteanu et al.; 2008a, 2009a).

Analiza comparativ-lingvistică a celor *patru* dicționare-tezaure a mai relevat faptul că fiecare dicționar are caracteristici particulare și necesită o modelare lingvistică și lexico-semantică atentă. Cu toate acestea, multe concepte referitoare la sensuri, definițiile sensurilor, marcherii de sensuri, exemplele la definiții, referirea sensurilor și siglele pot fi similare, adaptabile și transferabile între configurațiile SCD utilizate, depinzând și de poziția și / sau nivelul la care se situează acestea în cadrul unui arbore de sensuri al unei intrări.

În urma analizării rezultatelor parsării intrărilor din cele *patru* dicționare, am determinat că metoda de parsare propusă, prin folosirea a trei configurații SCD pentru parsarea intrărilor de dicționare-tezaur, este eficientă, robustă și portabilă, putând fi aplicată cu succes peste tezaure (semnificativ) diferite, cu rezultate foarte bune comparativ cu metodele tradiționale de parsare (standard Dictionary Entry Parsing).

Mulțumiri. Rezultatele modelării lexical-semantice și parsării dicționarelor-tezaure **DLR** și **DAR** din această lucrare au fost obținute în cadrul cercetărilor la grantul eDTLR – PNCDI 2, No. 91_013/18.09.2007.

Referințe bibliografice

- Comitetul de revizie a DLR (1952). Reguli de codificare pentru DLR. Institutul de Filologie „Iorgu Iordan”, Academia Română, București.
- Cristea, D., Răschip, M., Forăscu, C., Haja, G., Florescu, C., Aldea, B., Dănilă, E. (2007). *The Digital Form of the Thesaurus Dictionary of the Romanian Language*. In Proceedings of the 4th International IEEE Conference SpeD 2007.
- Curteanu, N., E. Amihăesei (2004). *Grammar-based Java Parsers for DEX and DTLR Romanian Dictionaries*. ECIT-2004 Conference, Iasi, Romania.
- Curteanu, N. (2006). *Local and Global Parsing with Functional (F)X-bar Theory and SCD Linguistic Strategy*. (I.+II.), Computer Science Journal of Moldova, Academy of Science of Moldova, Vol. 14 no. 1 (40):74-102; no. 2 (41):155-182.
- Curteanu, N., Trandabăț, D., Moruz, A. M. (2008a). *Extracting Sense Trees from the Romanian Thesaurus by Sense Segmentation & Dependency Parsing*, Proceedings of CogAlex Workshop, COLING 2008, pp. 55-63, ISBN 978-1-905593-56-9.
- Curteanu, N., Trandabăț, D., Moruz, A., Bolea, C., Husarciuc, M. (2008b). *Parsarea dicționarului eDTLR cu configurații SCD la arborele de sensuri și definiții*. Raport de cercetare la Grantul PNCDI 2, Nr. 91_013/18.09.2007, Faza 2008.
- Curteanu, N., Moruz, A., Trandabăț, D. Bolea, C., Spătaru, M., Husarciuc, M. (2009a). *Parsare arborilor de sensuri și segmentarea la definiții în dicționarul tezaur* (Ed. D. Trandabăț, D. Cristea, D. Tufiș) Resurse lingvistice și instrumente pentru prelucrarea limbii române, ConsILR-2008, Editura Univ. „Al. I. Cuza” Iași, p. 65-74.
- Curteanu, N., Trandabăț, D., Moruz, A., Bolea, C. (2009b). *Parsarea dicționarelor eDA și eDLR cu configurații SCD la arborele de sensuri și definiții*. Raport de cercetare la Grantul PNCDI 2, Nr. 91_013/18.09.2007, Faza 2009.
- Erjavec, T, Evans, R., Ide, N., Kilgarriff, A. (2000). The CONCEDE Model for Lexical Databases. Research Report on TEI-CONCEDE LDB Project, Univ. of Ljubljana, Slovenia.
- Hauser, R., Storrer, A. (1993). *Dictionary Entry Parsing Using the LexParse System*. Lexikographica 9 (1993), 174-219.
- Kammerer, M. (2000). *Wörterbuchparsing Grundsätzliche Überlegungen und ein Kurzbericht über praktische Erfahrungen*, <http://www.matthias-kammerer.de/content/WBParsing.pdf>
- Lemnitzer, L., Kunze, C. (2005). *Dictionary Entry Parsing*, ESSLLI 2005.
- Neff, M., Boguraev, B. (1989). *Dictionaries, Dictionary Grammars and Dictionary Entry Parsing*, Proc. of the 27th ACL Vancouver, British Columbia, Canada Pages: 91 – 101.
- Tufiș, D. (2001). From Machine Readable Dictionaries to Lexical Databases, RACAI, Romanian Academy, Bucharest, Romania.
- XCES TEI Standard, Variant P5 (2007): <http://www.tei-c.org/Guidelines/P5/>

REALIZAREA UNUI TREEBANK ROMÂNESC

CENEL-AUGUSTO PEREZ

Universitatea „Al.I.Cuza”, Facultatea de Litere, Iași – România

Universitatea „Al.I.Cuza”, Facultatea de Informatică, Iași – România;

cperez@info.uaic.ro

Rezumat

În lucrarea de față vom inventaria rezultatele achiziționării corpusului de sintaxă a limbii române, punând accent atât pe problemele întâmpinate pe parcursul achiziționării acestui corpus (se vor menționa și câțiva pași ai creării corpusului), cât și pe soluționările acestor probleme.

1. Introducere

Cercetătorii au ajuns la un consens, și anume, că un progres semnificativ și rapid poate exista atât în înțelegerea textului cât și în înțelegerea limbii vorbite prin investigarea acelor fenomene lingvistice ce apar foarte frecvent în fapte de limbă întâlnite în mod natural, fără constrângeri, și prin încercarea de a extrage automat informații despre limbă din corpusuri¹ foarte mari. Astfel de corpusuri încep să servească drept unelte importante de cercetare pentru oamenii de știință din domeniile procesării limbajului natural, a recunoașterii vorbirii și a sistemelor de înțelegere a limbajului vorbit, precum și în lingvistica teoretică.

Între resursele utilizate pentru studierea sintaxei limbilor naturale o componentă importantă sunt treebank-urile. Un treebank este un corpus anotat la sintaxă într-un formalism general acceptat, ca de exemplu, gramatici de dependență, bazate pe relații directe existente între cuvinte (Mel'cuk, 1987), adică „relații de la cuvânt la cuvânt între cuvinte individuale” (Hudson, 1998), unde toate cuvintele dintr-o frază, cu excepția unuia (care se numește rădăcina frazei) depind de alte cuvinte. Treebank-uri există pentru limbi precum: chineza, ceha, engleza, franceza, germana, italiana, japoneza, poloneza, portugheza, spaniola, turca etc. Cele mai cunoscute treebank-uri actuale sunt pentru limba engleză (Penn Treebank – construit la Universitatea din Pennsylvania, Philadelphia²) și pentru limba cehă (Prague Dependency Treebank – construit la Universitatea Charles din Praga³). În dezvoltarea unui parser sintactic pentru limba română (Seretan et al., 2010) s-a utilizat un treebank relativ limitat atât în ceea ce privește mărimea lui, cât și în ceea ce privește complexitatea structurii sintactice (adică nu conține propoziții subordonate, iar mărimea medie a propozițiilor este doar de 9 cuvinte). În aceeași situație se afla și treebank-ul lui Călăcean și Nivre (2000), unde întâlnim doar texte din articole de ziare (în special subiecte politice și administrative) cu

¹ Corpusuri (pl.) – corpus (sg.) = Un corpus este o colecție de fapte de limbă, sub formă de text electronic, selectate după criterii externe pentru a reprezenta, pe cât posibil, o limbă sau o varietate de limbă ca fiind o sursă de date pentru cercetare lingvistică. (Sinclair 2004)

² <http://www.cis.upenn.edu/~treebank/>

³ <http://ufal.mff.cuni.cz/pdt/>

o medie de 8,94 token-uri pe propoziție. În prezenta lucrare descriem aspecte legate de crearea unei corpus treebank (mai complex) pentru limba română.

În secțiunea a doua vom prezenta succint pașii creării corpusului, iar în secțiunile 3 și 4 sunt detaliate câteva dintre problemele (însoțite, bineînțeles, de rezolvări) tehnice, respectiv gramaticale, întâlnite în procesul de achiziționare a corpusului. Ultima secțiune o constituie concluziile unde se vor face o inventariere și o evaluare a rezultatelor.

2. Crearea corpusului treebank

Crearea unui corpus treebank presupune mai mulți pași. Primul pas al acestei construcții este achiziționarea surselor lexicale. Am ales texte ce reflectă orientarea spre o selecție a situațiilor sintactice tipice, din care putem să ne edificăm asupra esenței combinațiilor sintactice de diferite tipuri și pe care le putem utiliza ca modele în diverse analize sintactice mai mult sau mai puțin complicate. S-a urmărit ca exemplele să provină din diverse genuri literare și să surprindă particularitățile sintactice ale variatelor tehnici de creație care au caracterizat o epocă sau alta, ori, mai exact, care au caracterizat arta literară a diferiților scriitori. Multe din textele din care e constituit acest corpus sintactic (împreună cu o parte din textele FrameNet-ului englezesc, traduse în limba română) sunt preluate din manualele și cursuri practice de limba română, din scriitori consacrați pentru a asigura acuratețea gramaticală și diversitatea sintactică.

Astfel s-a făcut o selecție a unor texte, de regulă propoziții și fraze cu o întindere între 3 și 144 cuvinte din operele unor autori precum: Mihai Eminescu, Ion Creangă, Octavian Goga, Liviu Rebreanu, Marin Preda, Mihail Sadoveanu, Barbu Ștefănescu Delavrancea, Calistrat Hogaș, George Coșbuc, Camil Petrescu și alții. Aceste texte, împreună cu doar câteva citate din ziare și almanahuri, au fost fie tehnoredactate de către mine și transformate în format electronic (format *.doc*) cu ajutorul editorului Microsoft Word din cadrul pachetului de servicii Microsoft Office, fie scanate și transformate în text electronic prin intermediul unui program OCR (Optical Character Recognition⁴), cu efectuarea corectărilor de rigoare, deoarece operația de OCR-izare produce adesea greșeli de interpretare. Astfel, o serie de texte au fost făcute disponibile în format electronic.

Al doilea pas în crearea corpusului îl constituie marcarea automată a informațiilor de natură morfologică asupra corpusului cu ajutorul webservice-ului de etichetare morfo-sintactică RACAI⁵. De exemplu, luăm fraza: *Rolul covârșitor a lui Eminescu în dezvoltarea literaturii și a limbii noastre literare a determinat crearea, începând din acest an, a unei catedre speciale „Mihai Eminescu” la Universitatea din București.* Această frază, după procesarea textuală a webservice-ului va arata astfel:

RolullrollNSRYINemsry covârșitor|covârșitor|ASN|Afpms-n alallTS|Tsfsluillui|TSO|Tfso EminesculEminesculNP|Np în|în|S|Spsa dezvoltarealdezvoltare|NSRY|Ncfsry

⁴ O aplicație pentru calculator care transformă imagini scanate de texte în șiruri de coduri de caractere, ce vor putea fi astfel căutate, indexate și prelucrate.

⁵ <http://www.racai.ro/webservices/TextProcessing.aspx>

REALIZAREA UNUI TREEBANK ROMÂNESC

literaturii literatură NSOY Ncfsoy și și CRICrssp alal TS Tsfs limbii limbă NSOY Ncfsoy
 noastre mele PSSIDs lfsop literare literar ASON Afpfson alavea VA3SIVa--3s
 determinat determinat VPSM Vmp--sm crearea crearea NSRY Ncfsoy
 ,, ,, COMMA COMMA începînd începînd VG Vmg din din S Spsa
 acestacest DMSR Dd3msr---e anlan NSN Ncms-n ,, ,, COMMA COMMA alal TS Tsfs
 unei lun TSO Tifso catedre catedră NSON Ncfson speciale special ASON Afpfson
 ,, ,, DBLQ DBLQ Mihai_Eminescu Mihai_Eminescu NP NP ,, ,, DBLQ DBLQ
 lall S Spsa Universitatea universitate NSRY Ncfsoy din din S Spsa
 București București NP NP .l. PERIOD PERIOD

Transformat în xml, acest text va arăta astfel:

```

1 <?xml version='1.0' encoding='UTF-8'?>
2 <treebank id="Augusto">
3 <sentence id="1" parser="" user="Augusto" date="2010-30-19">
4 <word id="1" form="Rolul" lemma="rol" postag="Ncmsry" head="0" deprel=""/>
5 <word id="2" form="covârșitor" lemma="covârșitor" postag="Afpmn" head="0" deprel=""/>
6 <word id="3" form="a" lemma="a" postag="Tsfs" head="0" deprel="det."/>
7 <word id="4" form="lui" lemma="lui" postag="Tf-so" head="0" deprel=""/>
8 <word id="0" form="Eminescu" lemma="Eminescu" postag="Np" head="0" deprel=""/>
9 <word id="6" form="în" lemma="în" postag="Spsa" head="0" deprel=""/>
10 <word id="7" form="dezvoltarea" lemma="dezvoltare" postag="Ncfsoy" head="0" deprel=""/>
11 <word id="8" form="literaturii" lemma="literatură" postag="Ncfsoy" head="0" deprel=""/>
12 <word id="9" form="și" lemma="și" postag="Crssp" head="0" deprel=""/>
13 <word id="10" form="a" lemma="a" postag="Tsfs" head="0" deprel=""/>
14 <word id="11" form="limbii" lemma="limbă" postag="Ncfsoy" head="0" deprel=""/>
15 <word id="12" form="noastre" lemma="meu" postag="Dslfsop" head="0" deprel=""/>
16 <word id="13" form="literare" lemma="literar" postag="Afpfson" head="0" deprel=""/>
17 <word id="14" form="a" lemma="avea" postag="Va--3s" head="0" deprel=""/>
18 <word id="10" form="determinat" lemma="determina" postag="Vmp--sm" head="0"/>
19 <word id="16" form="crearea" lemma="creare" postag="Ncfsoy" head="0" deprel=""/>
    
```

Fișierul rezultat după această procesare a fost transformat în fișier *.xml* printr-un program în PERL, care face modificările necesare astfel încât fișierul să fie adaptat pentru adnotarea sintactică.

Deoarece utilizarea unui adnotator sintactic implică un proces îndelungat și câteodată anevoios, este necesară conceperea unei metodologii de adnotare la nivel sintactic a corpusului. Această metodologie reprezintă un alt pas în procesul creării.

Avem nevoie de instrucțiuni care să permită o adnotare la nivel sintactic consistentă cu o teorie lingvistică. Până în prezent s-a elaborat o listă cu posibii regenți, posibile cuvinte subordonate, relații etc. În această listă se mai poate completa pe măsură ce noi exemple sunt descoperite pe parcursul procesului de adnotare sintactică. De exemplu, în figura de mai jos (primul rând) avem un substantiv ca regent. S-a găsit o propoziție cu un grup nominal în care substantivul era urmat de o prepoziție și un adverb. În arborele de dependență, pe săgeata dintre substantiv (*plimbarea*) și prepoziție (*de*), vom avea notată, așadar, relația de atribut adverbial (*a.adv.*), ca în exemplul: *Plimbarea de azi*. Pentru continuarea analizei acestui grup nominal, trebuie să căutăm în listă situația în care regentul este o prepoziție (*de*) și cuvântul subordonat este adverb (*azi*) și să aplicăm aceeași metodologie pentru a completa arborele.

CENEL-AUGUSTO PEREZ

REGENT	CUV. SUBORDONAT	URMAT DE	RELATIE	ABREVIERE	EXEMPLU
1. Sub stantiv	prepozitie	adverb	atribut adverbial	a. adv.	Plimbarea de azi
2. Sub stantiv	prepozitie	verb nefinit	atribut verbal	a. verb.	Masina de spalat
3. Sub stantiv	prepozitie	nominal	atribut substantival	a. subst.	Praf de pușcă
4. Sub stantiv	numeral	-	atribut adjectival	a. adj.	Doi oameni
5. Sub stantiv	articol	-	determinant	det.	Un om
6. Sub stantiv	adjectiv	-	atribut adjectival	a. adj.	Fată tânără
7. Sub stantiv	pronume	-	atribut adjectival	a. adj.	Casa lui
8. Sub stantiv	verb	-	atribut verbal	a. vb.	Ideea de a merge
9. Adjectiv	prepoziție	-	complement de agent	c. ag.	Construit de Ion
10. Adjectiv	adverb	-	compl. circum. de mod	c. c. m.	construit repede

Figura 1: Fragment din lista relațiilor de dependență (în curs de dezvoltare)

După operațiile de preprocesare, se poate trece la cel din urmă pas și cel mai important pentru proiectul de față: adnotarea corpusului la nivel sintactic. Această adnotare manuală, dacă este lipsită de reprezentarea grafică a arborilor de dependență, devine un proces foarte dificil. Prin urmare, a fost nevoie de utilizarea unui instrument grafic interactiv cu care să se adnoteze, vizualizeze și modifice arborii rezultați în urma parsării.

TreeAnnotator este un program⁶ de editare de arbori de dependență funcțională care oferă utilizatorului o interfață grafică.

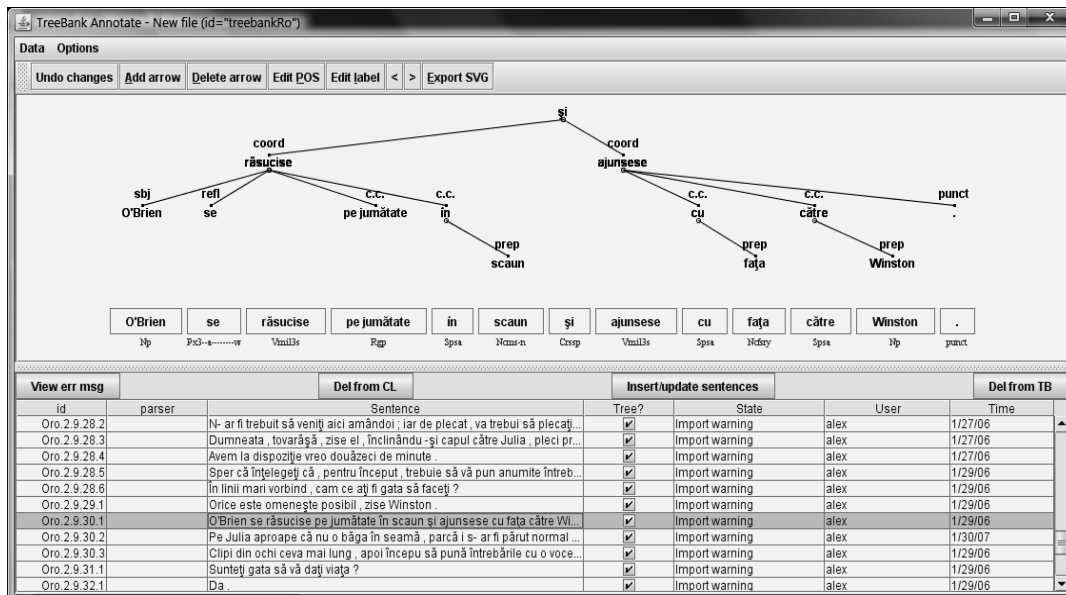


Figura 2: Arbore de dependență vizualizat cu TreeAnnotator

Interfața permite vizualizarea fiecărui cuvânt al frazei cu eticheta morfologică corespunzătoare. Această facilitate a fost introdusă pentru corectarea eventualelor probleme de adnotare a părții de vorbire. Pentru realizarea arborelui de dependență a

⁶ Acest program a fost dezvoltat în cadrul Institutului de Informatică Teoretică Iași și se găsește la adresa: <http://students.info.uaic.ro/~mmoruz/FDAnnotator/TreeAnnotator.tgz>.

unei fraze se selectează cuvântul părinte (regentul) urmat de cuvântul fiu (subordonatul), după care se inserează drumul de la părinte către fiu. În urma realizării tuturor legăturilor de dependență dintre cuvinte se obține arborele de dependență.

În continuare vom descrie diferitele tipuri de probleme ivite.

3. Probleme tehnice

S-a întâlnit problema segmentării greșite a unor fraze din cauza punctelor de suspensie sau a parantezelor pătrate cu puncte de suspensie. De exemplu, o frază ca:

„A cerceta evoluția scrierii chilirice românești, începînd din secolul al XVI-lea pînă spre mijlocul secolului al XIX-lea [...], înseamnă, evident, a-i urmări, și cronologic și pe regiuni istorice, mersul și dezvoltarea, precizînd momentele, direcția, semnele și proporția schimbărilor grafice, ca tot atîtea etape dialectice ale acestei forme de cultură, deosebit de importante.”⁷

a fost segmentată în două fraze, una conținînd textul: „A cerceta evoluția scrierii chilirice românești, începînd din secolul al XVI-lea pînă spre mijlocul secolului al XIX-lea [...]” și cealaltă avînd textul: „înseamnă, evident, a-i urmări, și cronologic și pe regiuni istorice, mersul și dezvoltarea, precizînd momentele, direcția, semnele și proporția schimbărilor grafice, ca tot atîtea etape dialectice ale acestei forme de cultură, deosebit de importante.” În realitate cele două întinderi constituie o frază doar luată împreună.

Soluția găsită acestei probleme a fost să modificăm textele inițiale și să scoatem punctele de suspensie sau parantezele pătrate cu puncte de suspensie din structura originală a frazei, deoarece acestea nu schimbă cu nimic analiza sintactică a textului, și apoi să trecem frazele prin marcarea automată cu informații lingvistice.

Similar s-au rezolvat și problemele: texte fără punct la sfîrșitul propoziției / frazei, ghilimele, propoziții / fraze ce încep cu literă mică.

Deoarece programul adnotator, TreeAnnotator, nu permite adnotarea de noduri izolate, fără legături cu restul nodurilor din frază, iar orice semn de punctuație constituie, la rîndul lui, un nod, s-a pus problema relaționării punctului de la sfîrșitul frazei / propoziției. De cine trebuie el legat? Am hotărât să legăm punctul de verbul predicativ al ultimei propoziții principale, deoarece un nod nu poate avea mai mulți părinți.

Ne-am mai lovit, totodată, de câteva probleme de procesare (unele texte produceau erori pe parcursul rulării diferitelor programe) și de timp (unele procesări au consumat resurse de timp semnificative), însă în capitolul următor vom descrie unele probleme de natură gramaticală.

4. Probleme gramaticale

Un rol important în analiza gramaticală îl joacă flexiunea, iar prin flexiune ne referim la totalitatea schimbărilor suferite de un cuvânt în aspectul lui formal pentru exprimarea raporturilor sintactice din vorbire între obiectele (ființe și lucruri) din realitatea

⁷ Text preluat din N. Angheliescu Temelie, Silviu Constantinescu – *Analize gramaticale*, Ed. Științifică și Enciclopedică, București, 1985, p.10-99;

înconjurătoare, precum și între obiecte și acțiuni. Prin urmare, substantivul sau substitutul său (numele obiectului) și verbul (cuvântul care exprimă acțiunea, starea sau existența privite ca proces) își vor modifica forma pentru a putea exprima raporturile sintactice care reflectă raporturile existente în realitatea înconjurătoare. Analiza trebuie să se facă deci, în strânsă legătură cu funcția pe care o îndeplinesc în propoziție aceste părți de vorbire. Plecând de la importanța funcției și de la faptul că un nod nu poate avea mai mulți părinți, am întâmpinat problema elementului predicativ suplimentar. Dacă luăm ca exemplu propoziția „*Azi, ca un sfânt dintr-o icoană veche,*

Blînd îmi răsai cu fața ta blajină”, cuvântul *blînd*, fiind considerat element predicativ suplimentar, ar putea avea ca regent și verbul *răsai* (cum răsai? blând) și substantivul *sfînt* (ce fel de sfânt? blând). Astfel că, am hotărât să legăm cuvântul respectiv doar de verb și să-i atribuim funcția sintactică de complement circumstanțial de mod.

O altă problemă este cea a conjuncțiilor subordonatoare fără funcție sintactică. Ele nefiind analizate gramatical în analizele tradiționale ale unor specialiști, s-a pus problema de cine le legăm și ce relație le atribuim. Deoarece conjuncția respectivă introduce o propoziție, ni se pare firesc să fie legată de verbul predicativ al propoziției respective, iar relația de pe arcul dintre verb și conjuncție să fie cea de particulă, conjuncția fiind considerată, aici, o particulă a verbului respectiv. În ceea ce privește conjuncțiile coordonatoare, acestea pot fi noduri părinte (chiar la primul nivel) în arborele sintactic din adnotarea propriu-zisă, fiindcă ele leagă două propoziții de același fel, dar nu țin numai de una sau de cealaltă.

Cea mai des întâlnită problemă a fost cea a locuțiunilor de orice fel. Programele folosite la prepararea textelor (vorbim despre procesele de pre-adnotare), în foarte multe cazuri, nu recunoșteau locuțiunile ca pe un singur nod, ci constituia câte un singur nod pentru fiecare cuvânt inclus în locuțiune. De exemplu, *pas cu pas*, o locuțiune adverbială, este recunoscută de program în trei noduri separate *pas*, *cu* și *pas*, ceea ce pune în dificultate analiza noastră sintactică, care de altfel ar fi fost foarte simplă dacă am considera locuțiunea ca fiind un singur nod. Pentru a simplifica lucrurile, am umblat, din nou, de data asta, la textul din *.xml* și am legat cuvintele din locuțiune prin caracterul „underscore” („_”) pentru ca programul să recunoască șirul de cuvinte *pas_cu_pas* ca pe un singur nod. Bineînțeles a trebuit să modificăm și denumirea părții de vorbire, ea fiind acum, în cazul de față, cea de adverb.

5. Concluzii

În concluzii, după ce am trecut în revistă unele probleme și soluțiile lor, nu ne rămâne decât să tragem linia și să punem pe masă ceea ce s-a realizat până în prezent. S-a pornit de la o colecție de aproximativ 2000 de texte (fraze și propoziții), însumând în jur de 70.000 de cuvinte și s-a ajuns la un număr de 15 fișiere *.xml*, care conțin peste 1.800 de fraze adnotate cu aplicația TreeAnnotator, având un total de aproximativ 67.000 de cuvinte.

Această colecție de texte adnotate constituie doar un început, un punct de pornire pentru a demara învățarea, la rândul nostru, a unui parser sintactic pentru limba română. Vom continua cu adnotările la nivel sintactic, dar de această dată automat cu ajutorul aplicației MALT (Nivre et al., 2007), care creează un model de adnotare la dependență pentru limba română, până când se va atinge un număr semnificativ de cuvinte,

deoarece, cu cât programul se va antrena pe un număr cât mai mare de texte, cu atât el va avea o bază de date mai bogată și rezultatele, în urma unor teste, vor fi mai eficiente.

Un corpus adecvat, de mai multe milioane de cuvinte în uz este dincolo de puterea de organizare a unui singur cercetător, el necesitând o activitate în cadrul unui colectiv. De aceea, corpusul, realizat până în acest stadiu din cercetarea de față, constituie un „start” promițător pentru ceea ce va urma în viitorul apropiat, și anume un corpus *gold*. Această denumire o poartă doar acele corpusuri ce se supun unor condiții de realizare, și anume: corpusurile trebuie să conțină un număr foarte mare de cuvinte (de ordinul milioanei de cuvinte) care să fie adnotate de cel puțin doi adnotatori umani (bineînțeles, aceștia trebuie să fie specialiști în domeniu).

Referințe bibliografice

A. Izvoare și lucrări de referință

- Angheliescu Temelie, N., Constantinescu, S. (1985). *Analize gramaticale*, Editura Științifică și Enciclopedică, București, p.10-99.
- Angheliescu Temelie, N., Popescu, A. G. (1976). *Dificultăți ale analizei gramaticale*, Editura Științifică și Enciclopedică, București, p.131-291.
- Botiș, V., Vulișici Alexandrescu, M., Comănescu, I. (1977). *Sintaxa propoziției*, Editura Facla, p. 50-103.
- Constantinescu, S. (1976). *Exerciții și analize gramaticale*, Editura Didactică și Pedagogică, București, p.61-80.
- Drașoveanu, D. D. (1959-1966). *Analize gramaticale și stilistice*, Editura Științifică, București, p. 111-155.
- Lupu, C., Vlădescu, A. (1997). *Limba română prin exerciții*, Editura Logos, București, p. 18.
- Metea, A., Drincu, S. (1974). *Analize gramaticale*, Editura Facla, Timișoara, p.23-109.
- Muțiu, I., Bercea, L. P. (1985). *Exerciții de gramatică*, Editura Facla, Timișoara, p.33-77.
- Vlad, V., Știrbu, P., Vlad-Budoiu, V. (1978). *Analize gramaticale structurale*, Editura Dacia, Cluj-Napoca, p.57-550.

B. Literatură secundară

- Academia Română (2005). *Gramatica limbii române*, Editura Academiei Române, București.
- Călăcean, M. and Nivre, J. (2009). A Data-Driven Dependency Parser for Romanian. In Proceedings of TLT-7.
- Curteanu, N., Gâlea, D., Butnariu, C. (2003). *Segmentation Algorithms for Clause-Type Textual Units*, versiunea în limba română în „Limba română în societatea informațională”, D. Tufiș, Fl. Filip Eds., Romanian Academy, Bucharest, Expert Press, p.165-190.
- Curteanu, N., Moruz, A., Trandabăț, D., Bolea, C., Dornescu, I. (2006). The Structure and Parsing of Romanian Verbal Group and Predicate, *Advances in Intelligent*

- Systems and Technologies ECIT2006 – 4th European Conference on Intelligent Systems and Technologies, Iasi, Romania, Septembrie 21-23, pp. 93-105.
- Hudson, R., (1998). English Grammar. London, Rontledge.
- Mel’Cuk, I. A., (1987). Dependency Syntax: Theory and Practice. Buffalo, Suny Press.
- Nivre, J. (2007). Data-Driven Dependency Parsing across Languages and Domains: Perspectives from the CoNLL 2007 Shared Task. In Proceedings of the Tenth International Conference on Parsing Technologies, 168-170.
- Seretan, V., Wehrli E., Nerima L., Soare G. (2010). (submitted for publication) FipsRomanian : Towards a Romanian Version of the Fips Syntactic Parser. LREC 2010.
- Stati, S. (1972). *Elemente de analiză sintactică*, Editura Didactică și Pedagogică, București.
- Tufiș, D., Barbu, A. M., Pătrașcu, V., Rotariu, G., Popescu, C. (1997). *Corpora and Corpus-Based Morpho-Lexical Processing*, in „Recent Advances in Romanian Language Technology”, Dan Tufiș, Poul Andersen, Editura Academiei Române.
- Tufiș, D. (2000). *Using a large set of EAGLES-compliant Morpho-Syntactic Descriptors as a Tagset for Probabilistic Tagging*, Proceedings of LREC.
- <http://www.di.unito.it/~tutreeb/>
- <http://consilr.info.uaic.ro/ro/index.php?showpage=030102>
- <http://www.ceid.upatras.gr/Balkanet/publications.htm>
- <http://www.sketchengine.co.uk/>
- <http://ufal.mff.cuni.cz/pdt2.0/>
- http://www.acad.ro/pro_pri/doc/st_h01.doc
- <http://students.info.uaic.ro/~mmoruz/FDAnnotator/TreeAnnotator.tgz>
- <http://www.racai.ro/webservices/TextProcessing.aspx>

MONITORIZAREA PRESEI ÎN CADRUL PROIECTULUI NEOROM

ANA-MARIA BARBU

Institutul de Lingvistică „Iorgu Iordan – Al. Rosetti”, București – România
anamaria.barbu@g.unibuc.ro

Rezumat

Acest articol prezintă contribuția părții române în cadrul proiectului Neorom. Proiectul are ca obiectiv dezvoltarea unei baze de date cu neologismele lexicale adunate din ziare, începând cu anul 2004, în toate limbile romanice. Articolul are trei părți. În prima parte este prezentat sistemul românesc de monitorizare a presei, în a doua, interfața Neorom de înregistrare a neologismelor în baza de date, iar în a treia, o scurtă analiză a cuvintelor din bază din perspectiva modului lor de formare.

1. Introducere

Proiectul Neorom (proiect de Neologie romanică) este un proiect internațional care are ca obiectiv crearea unei baze de date multilingve în care sunt înregistrate cuvintele nou apărute în presa scrisă, în paralel pentru următoarele limbi romanice: catalană, franceză (din Franța, Belgia și Québec), galiciană, italiană, portugheză (din Portugalia și Brazilia), română și spaniolă. Inventarierea cuvintelor a început în anul 2004 și continuă în fiecare an.

Acest proiect se desfășoară sub egida organizației Realiter (<http://www.realiter.net>) și a Uniunii Latine (<http://www.unilat.org>) și este condus de Maria Teresa Cabré, profesor la Institutul pentru Lingvistică Aplicată din cadrul Universității Pompeu Fabra, Spania. El este concretizat printr-o platformă de căutare a cuvintelor inventariate în bazele de date corespunzătoare limbilor implicate (Cabré 2004). Platforma Neorom se găsește la adresa <http://obneo.iula.upf.edu/bneorom/index.php>.

Proiectul Neorom este unul dintre proiectele grupului L'Observatori de Neologia (<http://www.iula.upf.edu/obneo/obpresca.htm>) cu cea mai importantă activitate în domeniul neologiei și al terminologiei. Mai putem menționa ca proiecte similare proiectul Neoscope (<http://www.certa.usj.edu.lb/files/neoscope.htm>) și secțiunea de achiziție de cuvinte noi *Oxford English Dictionary*, <http://dictionary.oed.com/about/>.

În accepția proiectului Neorom, cuvintele nou apărute, numite și *neologisme*¹, sunt acele unități lexicale formate din unul sau mai multe cuvinte (grafice) care nu sunt înregistrate într-un corpus prestabilit de dicționare reprezentative ale limbii considerate, numit *corpus de excludere*. Prin urmare, cuvintele inventariate în baza românească Neorom nu sunt neapărat cele care apar în presă strict într-un anumit an, ci și acelea care dintr-un motiv sau altul nu au fost înregistrate în dicționare, deși sunt folosite în limbă de mai

¹ Termenul de *neologism* nu are o semnificație unanim acceptată, de aceea pentru unii lingviști folosirea lui în contextul acestui proiect poate părea improprie. O dezbateră asupra semnificației acestui termen ar depăși spațiul acestui articol.

multă vreme. Pe de altă parte, proiectul acordă deopotrivă atenție cuvintelor care nu au intrat propriu-zis în limbă și e posibil să nu intre niciodată, dar reprezintă mărturie asupra mijloacelor de îmbogățire a limbii și asupra inventivității lingvistice. Altfel spus, nu se aplică o selecție a cuvintelor în funcție de numărul lor de apariții în corpusul de ziare considerat, ci se inventariază până și hapaxurile. De asemenea, trebuie precizat că în baza de date Neorom înregistrăm și cuvinte care par să existe în dicționare, dar care au suferit totuși o modificare; de pildă, și-au schimbat partea de vorbire sau genul, au căpătat forme de plural sau folosesc alte morfeme de plural, se scriu altfel (cu cratimă sau fără cratimă sau ca în limba de origine sau reprezintă greșeli de ortografie) ș.a.

Cercetarea prezentată în acest articol constă în dezvoltarea bazei de date românești din cadrul proiectului Neorom. În secțiunea următoare descriem sistemul de monitorizare folosit pentru inventarierea cuvintelor ce urmează a fi introduse în bază. În secțiunea 3 prezentăm interfața bazei de date Neorom cu câmpurile din fișa care trebuie completată pentru fiecare cuvânt introdus. Unele dintre câmpuri constituie criterii de căutare în bază, cunoașterea lor fiind necesară celor care vor dori să utilizeze baza Neorom. Secțiunea 4 o dedicăm unei scurte analize a rezultatelor românești din cadrul proiectului. Utilitatea acestui proiect și în special a bazei de date românești este subliniată în capitolul de concluzii.

2. Sistemul de monitorizare

a. Date preliminare

Corpusul folosit pentru monitorizarea cuvintelor nou-apărute în presă este alcătuit din mii de articole (însușind circa 14 milioane de cuvinte) culese de pe Internet din variantele on-line a mai multor ziare precum *Adevărulonline*, *BBC-România*, *Cotidianul*, *Jurnalul Național*, *Libertatea*, *România Liberă* și altele, apărute în perioada 01.08.2004 – 15.10.2007. Corpusul este accesibil prin bunăvoința unor colegi italieni. Trebuie subliniat faptul că obținerea textelor dintr-un ziar on-line, într-un format care poate fi prelucrat automat ulterior, nu este deloc un lucru trivial, de aceea ne-am bucurat să putem beneficia de ajutorul colegilor noștri. Începând cu anul acesta vom încerca să folosim programul Buscaneo, pus la dispoziție de colegii spanioli ai acestui proiect, însă avem reticente în privința eficienței lui. Pentru fiecare cuvânt din ziarul inspectat pe care programul nu-l găsește în dicționarul de forme flexionare, Buscaneo deschide o fișă Neorom. Dicționarul nostru electronic fiind relativ mic, numărul de cuvinte negăsite este foarte mare, prin urmare și numărul de fișe deschise va fi covârșitor și inutil de mare, deoarece vom fi privați de eficiența lucrului cu liste de cuvinte pe care le putem ulterior prelucra automat sau le putem inspecta vizual.

Este de menționat că *BBC-România* folosește diacritice, celelalte ziare nu folosesc diacritice. Toate ziarurile folosesc grafia cu „â” (care în textul fără diacritice apare „a”), cu excepția ziarului *Cotidianul* care folosește și grafia cu „î” (care în textul fără diacritice apare „i”).

Dicționarele românești care alcătuiesc corpusul de excludere, adică referința față de care un cuvânt este considerat nou apărut (anume, dacă nu se găsește deja înregistrat în corpusul de excludere respectiv) sunt cele în format electronic disponibile la adresa

<http://dexonline.ro>, din care avem în vedere DEX '98, DLRA, DLRC, DLRM, DN, MDA și MDN, precum și DOOM2, pentru că prezintă avantajul căutării semiautomate.

b. Procesul de monitorizare

Modul în care selectăm cuvintele nou apărute este unul relativ simplu, după cum se poate vedea în figura 1, însă destul de laborios prin volumul mare de date și prin combinarea etapelor automate, semiautomate și manuale. Articolele din ziare sunt indexate cu ajutorul concordanțierului LUCON (<http://sourceforge.net/projects/lucon/>). Cuvintele (flexionate) din indexul obținut sunt căutate automat în dicționarul electronic de forme flexionare RoMorphoDict (Barbu, 2008), iar cele care nu au fost găsite în acest dicționar (care conține paradigmele complete ale celor aprox. 60000 de cuvinte-titlu înregistrate în DOOM 1) sunt trecute într-o listă de cuvinte „necunoscute”. Această listă poate fi impresionant de mare, de ordinul zecilor de mii de cuvinte. O formă necunoscută poate fi înregistrată astfel din următoarele motive:

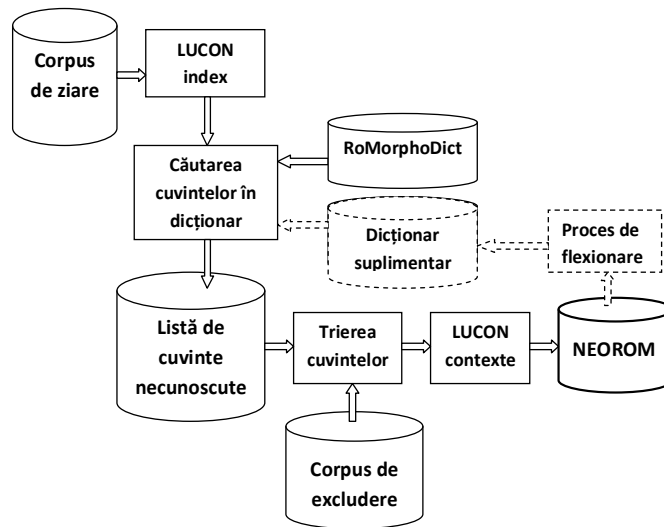


Figura 1: Sistemul de monitorizare

- forma respectivă nu se află în dicționarul electronic folosit ca referință (pentru că acesta conține doar 60000 de leme, multe dintre acestea nefiind folosite în limbajul curent);
- face parte dintre abrevieri, unități de măsură sau alte simboluri;
- este un nume propriu;
- este un cuvânt străin;
- reprezintă o greșeală tipografică sau a vorbitorilor (de pildă, forma *ostatec* este întâlnită de 111 ori);
- exprimă perioade și aproximări (de exemplu, ianuarie-aprilie, primăvară-vară, opt-nouă).

Această listă de cuvinte „necunoscute” este apoi inspectată vizual în etapa de triere a cuvintelor, pentru a exclude acele cuvinte pe care le socotim reziduale, iar ceea ce rămâne este confruntat, de asemenea, manual, cu dicționarele amintite din corpusul de

excludere, care au totuși avantajul că sunt introduse într-o bază de date accesibilă pe Internet. Pentru fiecare cuvânt care nu a fost găsit în corpusul de excludere se completează o fișă de introducere în baza de date Neorom, după cum e prezentat în secțiunea următoare.

Procesul descris până aici trebuie completat cu o nouă etapă în care cuvintele care au fost introduse în baza Neorom sunt flexionate și înregistrate într-un dicționar suplimentar de referință, pentru a nu mai fi inventariate încă o dată. De fapt întregul proces poate fi îmbunătățit dacă am avea la dispoziție un program de descărcare a textelor de ziare de pe Internet, dacă corpusul de excludere ar fi în formă electronică corespunzătoare sau dacă am dispune de un program de flexionare a cuvintelor necunoscute pe care să-l putem integra într-un proces de prelucrare.

3. Interfața platformei Neorom

În această secțiune descriem conținutul fișei care trebuie completată pentru fiecare cuvânt nou introdus în bază. După cum se vede în figura 1, câmpurile care trebuie completate sunt următoarele.

Date generale despre **despuiere**:

Tipo – Se indică faptul că datele sunt extrase din presa scrisă.

Lengua – Datele culese sunt din limba română.

Ámbito geográfico – Este vorba de limba română vorbită în România.

Fuente – Ziarul din care s-a extras cuvântul respectiv este *Libertatea* (varianta electronică).

F. Publicación – Data ediției ziarului este 06 ianuarie 2005.

Datele fișei:

Entrada – Se indică cuvântul-titlu.

Categoría gramatical – Categoria gramaticală (partea de vorbire, genul, numărul, tipul verbului etc.) se alege dintr-o listă predefinită. În privința categoriei gramaticale, ne-am confruntat la început cu faptul că în lista predefinită nu exista genul neutru (deoarece nu există în celelalte limbi romanice), însă ulterior această problemă a fost rezolvată.

Contexto – Pentru fiecare cuvânt-titlu se dă contextul în care apare în numărul respectiv al ziarului. În cele mai multe cazuri, un cuvânt apare în mai multe numere ale mai multor ziare. În această situație se alege de obicei numărul cel mai vechi, însă acolo unde este cazul se preferă contextul care sugerează cel mai bine sensul cuvântului respectiv (de pildă, sunt contexte în care se dă chiar o scurtă definiție a cuvântului respectiv). Contextele sunt obținute cu ajutorul concordanțerului Lucon, care, pe lângă multe alte facilități, o are pe aceea că oferă o cale directă și simplă de obținere a sursei contextului, adică ziarul și ziua ediției (date ce trebuie completate în fișă).

Aspectos tipográficos – Această rubrică se referă la modalitățile de subliniere a unui cuvânt. În cele mai multe cazuri, sublinierea se face cu ghilimele sau paranteze, însă în edițiile pe hârtie pot exista și alte modalități.

Tipo de neologismo – Acest câmp este destinat indicării modului de formare a cuvântului nou. Se alege dintr-o listă predefinită simbolul corespunzător modului de formare a cuvântului. Conținutul listei predefinite, precum și o analiză detaliată a rezultatelor sunt prezentate în (Barbu 2010).

MONITORIZAREA PRESEI ÎN CADRUL PROIECTULUI NEOROM

PLATAFORMA NEOROM Francès · Galleg · Portuguès

Bienvenido/a ivr | **Gestión de fichas** :: Prensa - Editar ficha

- **Gestión de fichas**
- Edición de fichas
- Consulta realizadas
- Consulta general
- Consulta impresión
- **Salir**

Datos generales del vaciado

Tipo: Prensa

Lengua: Rumano

Ámbito geográfico: ROMANIA

Fuente: Libertatea

F. Publicación: 06/01/2005 (dd/mm/aaaa)

Datos de la ficha

Entrada: designeră corpus de exclusión

Categoría gramatical: Nombre femenino

Contexto: Valentina Pelinel: La o prezentare, *designera* nu mi-a inchis bine gacile la fusta si am ramas fara ea.

Aspectos tipográficos: Sin marca tipográfica

Tipo de neologismo: SINT - Neologismo sintáctico ayuda estadísticas

Sección del periódico: CULT - culturelle

Página:

Autor: Sí No

Código nota: 2g buscar

Nota: variante sintáctica de s.m. „designer”

Estado: Incompleta Completa

Comentarios:

Historial de cambios

	Creación	Completado	Corrección	Validación
Autor	ivr	-	-	-
Fecha	20/01/2009	-	-	-

Observatori de Neologia - Institut Universitari de Lingüística Aplicada - Universitat Pompeu Fabra
Roc Boronat, 138, 08018 Barcelona. Tel. 935 422 322 Fax 935 422 321
observatori.neologia@upf.edu

Figura 2: Fișa unei intrări din baza de date Neorom

Sección del periódico – Se indică secțiunea articolului din care a fost extras cuvântul. Domeniile cărora le aparțin în general secțiunile unui ziar sunt cel politic, economic, financiar, social, comercial de autoturisme, economic, informatic-electronic, sănătate, sport, cultural și de diverse. Pentru că în corpusul electronic folosit articolele nu sunt organizate pe secțiuni, domeniul este dedus după conținutul articolului.

Página – În cazul edițiilor pe hârtie se poate indica pagina publicației, însă noi, folosind edițiile electronice, nu avem acces la această informație.

Autor – Dacă se cunoaște autorul articolului sau cine a utilizat cuvântul (în cazul unui citat), acest fapt poate fi indicat.

Código nota – Pentru anumite tipuri de neologisme se poate adăuga un cod de notă explicativă, ales dintr-o listă predefinită.

Nota – Acest câmp conține nota explicativă propriu-zisă. În exemplul dat în figura 1, cuvântul *designeră* este un neologism sintactic (pentru definiție vezi secțiunea următoare), pentru că este obținut prin schimbarea genului de la substantivul masculin *designer* la unul feminin.

Estado – Starea fișei poate fi precizată în această rubrică prin indicarea uneia dintre valorile Incompletă/Completă.

Comentarios – Colaboratorii, care sunt membri ai proiectului, pot face diferite comentarii pe marginea fișei respective, ele nefiind făcute publice.

Historial de cambios – Zona de evidență a modificărilor operate în fișă este completată automat de sistem cu numele colaboratorului care a creat, a completat, a corectat sau a validat fișa și cu data operației respective.

Procesul de realizare a fișei se încheie apăsând butonul de salvare 'Guardar' sau lăsând fișa neschimbată prin butonul de anulare 'Cancelar'. Ștergerea fișei din baza de date se face prin butonul 'Borrar'.

4. Scurtă analiză a rezultatelor

Până în momentul de față au fost introduse în baza Neorom circa 2000 de cuvinte, urmând ca până la sfârșitul anului 2010 inventarul să fie extins la peste 4000. În faza actuală, completarea bazei de date se poate face într-un ritm de 100 de intrări pe zi. Clasificarea modului de formare a cuvintelor a trebuit să-l facem după standardul propus în cadrul proiectului, comun pentru toate limbile implicate. Pentru a ne alinia la acest standard, a fost necesar să facem unele compromisuri prezentate în cele ce urmează.

- Clasificarea cuvintelor după modul de formare se bazează pe forma lor analizabilă, iar nu pe etimologia lor în mod strict (de pildă, unele dintre cuvintele analizabile pot fi de fapt împrumuturi).

- Dintre mai multe moduri de formare posibile s-a ales unul singur. În principiu, s-a ținut seama de cuvintele deja înregistrate în dicționare ca bază de derivare. De exemplu, cuvântul nou *coprezentatoare* a fost considerat ca reprezentând derivarea cu prefix a lui *prezentatoare*, mai degrabă decât sufixarea moțională de la *coprezentator*, care nu e nici el înregistrat. Alt exemplu: un cuvânt precum *sepepist* e format atât prin siglare, cât și prin sufixare, dar s-a înregistrat în baza Neorom doar format prin siglare.

- Anumite moduri de formare pot diferi de ceea ce propune literatura de specialitate românească. De pildă, în proiectul Neorom femininul *designeră* este considerat un neologism sintactic (sau gramatical), în timp ce în literatura românească el este un derivat prin sufixare. Un alt exemplu este acesta: dacă înlăturând fie prefixul, fie sufixul aceluiași cuvânt obținem tot câte un neologism atunci cuvântul este considerat un derivat prin 'sufixare și prefixare' în Neorom (ex. *superblatist*: atât *superblat*, cât și *blatist* sunt neologice), pe când în literatura românească asemenea derivate nu sunt decât cele în care sufixul și prefixul sunt atașate unei baze inexistente ca atare în limbă (ex. *îm+brăc+a*).

- Proiectul Neorom consideră că prefixoide care în limba de origine au fost prepoziții sau adverbe (ex. *anti-*, *extra-*, *intra-*) formează derivate prin prefixare, iar celelalte prefixoide sunt elemente de compunere cultă (ex. *auto-*, *agro-*, *bio-*).

- Din păcate, proiectul Neorom nu oferă posibilitatea notării cuvintelor formate prin substituție de sufixe (ex. *atoputere* < *atoputernic*) sau cele prin derivare regresivă (*bâlbă* < *bâlbâi*).

Pentru stabilirea modurilor de formare, pe lângă recomandările din cadrul proiectului, am folosit și lucrările (Coteanu et al., 1985) și (Stoichițoiu-Ichim, 2001). Precizăm însă că modul de formare indicat pentru fiecare cuvânt din bază nu constituie rezultatul unei cercetări riguroase de etimologie, fapt care lasă loc de studiu specialiștilor interesați.

În cele ce urmează este prezentată succint repartizarea cuvintelor din baza de date românească din acest moment, în funcție de modurile de formare cele mai reprezentative (pentru care indicăm simbolurile folosite în platforma Neorom).

Derivare cu afixe: 27,5%

Derivare cu sufix – FSUF: 13,54%. Cele mai frecvente sufixe sunt **-ist** (ex. *cablist* „lucrător la cablu”, *galerist* „proprietar de galerie de artă”, *forumist* „utilizator de forumuri”), **-re** (ex. *fumare*, *gudurare*, *microcipare*) și **-(e/i)an** (ex. *becalian*, *cașavencian*). Cel mai neobișnuit sufix este **-ing** (ex. *bancheting* „organizare de banchete”, *cuponing* „instrument de marketing ce presupune folosirea cupoanelor”).

Derivare cu prefix – FPRE: 13,96%. De departe cel mai des utilizat prefix este **ne-** (30% din derivatele cu prefix) (ex. *nealegere*, *necalitativ*, *nesustenabil*) urmat îndeaproape de prefixoidul prepozițional **anti-** (22,5%), care este frecvent folosit pentru a forma adjective din substantive (ex. *anticutremur*, *antisărăcie*, *antisforăit*). Un alt prefix îndrăgit de vorbitori este **cvasi-** (ex. *cvasi-absență*, *cvasinațional*, *cvasi-război*).

Compunere: 27,9%

Compunerea din cuvinte întregi – FCOM: 17%. Acest mod este folosit cu precădere în exprimarea relațiilor politice, economice sau sociale (ex. *americano-canadian*, *chinois-indian*, *etnico-religios*), a relațiilor economice sau organizatorice (ex. *administrativ-teritorial*, *cerere-ofertă*, *financiar-contabil*), în exprimarea numelor de funcții sau meserii (ex. *director-coordonator*, *cazangiu-tinichigiu*) și altele.

Compunerea cultă (sau savantă) – FCULT: 8,8%. Acest tip de compunere se face cu ajutorul prefixoidelor (altele decât foste prepoziții sau adverbe) dintre care, în mod categoric, cel mai bine reprezentat este prefixoidul neologic **euro-** (ex. *euroasistat*, *eurobarometru*, *euromătură*, *europubelă*). Dintre prefixoidele clasice este folosit cu precădere **auto-**, cu sensul „de la sine, pe sine” (ex. *autoadresa*, *autobronzant*, *autoironic*). De remarcat este, de asemenea, apariția elementelor noi, de origine engleză, **tech-** (ex. *tech-cultură*, *tech demo*) și **web-** (*web designer/web-designer*, *web hosting*, *website*).

Compunerea prin abreviere (siglare – FTSIG, acronimie – FTACR, trunchiere – FTABR): 2,1%. Acest fel de compunere înregistrează două manifestări principale. Una este aceea de combinare a siglării numelor de partide sau organizații cu sufixul **-ist** pentru desemnarea membrilor grupării respective: *cederist*, *cekist*, *sepepist/sepepistă*, *udemerist*. Foarte rar se găsesc cuvinte formate prin siglare pură ca *sereleu* sau *sms*. A

doua manifestare privește trunchierea sau decuparea prefixoidelor din compusele culte pentru a fi folosite ca adjective invariabile sau, apoi, ca substantive. Exemple ale aplicării acestui procedeu înregistrate în baza de date sunt *combo, eco, ego, electro, etno, hidro, homo, lesbi, macro, mix, termo, turbo, vice*.

Conversiune – FCONV: 4,4%

Conversiunea este procedeul prin care un cuvânt este format prin schimbarea părții de vorbire, fără modificarea bazei lexicale. Substantivizarea adjectivelor este, de departe, tipul de conversiune cel mai des întâlnit. Se pot constata câteva direcții principale de manifestare a acestui tip de conversiune.

- când o persoană (sau, mai general, o ființă) este desemnată printr-o calitate a ei, mai ales la plural: *anchetat, mascat, național, necăjit, penal, vestic* etc.
- când persoanele sunt desemnate prin culorile care le definesc ca grup, fie că e vorba de o culoare politică: *verzii, alb-albaștri*, fie că e vorba de culorile unei echipe: *alb-violet, roșu-albaștri, tricolori, vișinii* etc.
- când se urmărește denumirea unei trăsături prin ceea ce are generic, caz în care adjectivul care denumește trăsătura respectivă este convertit în substantiv la singular, articulat hotărât: *derizoriul, divinul, economicul, eroicul, eroticul, esențialul, eticul*.
- când este suprimat elementul determinat de adjectiv dintr-o sintagmă uzuală: *(alegeri) anticipate, generale, legislative, locale; (aparate) electronice, (medicamente) compensate; (scor) egal* etc.

Neologisme sintactice (sau gramaticale) – SINT: 2.3%

Neologismele sintactice sunt cele care implică o schimbare de subcategorie gramaticală (gen, număr, regim verbal etc.) într-o bază lexicală. Acestea vizează în special substantivele masculine care capătă forme pentru genul feminin (ex. *administratoră, arbitură, doctoră, fană, stelistă*), substantivele neutre care devin masculine prin forma de plural (ex. *amortizori, boscheți, flotori*) sau substantivele care primesc forme de plural (ex. *bejuri, designuri, disconforturi, euroi*).

Variații ortografice sau flexionare – FVAR: 9,3%

Sunt socotite variații ortografice în special cuvintele care sunt scrise altfel decât în forma normată sau înregistrată în corpusul de excludere: *aiatolah* (normat, *ayatolah*), *afrodisiac* (*afrodiziac*), *chic* (*șic*), *click* (*clic*), *boutique* (*butic*). De asemenea s-au înregistrat numeroase cuvinte în care prefixul sau prefixoidul este despărțit prin cratimă de bază când norma recomandă scrierea legată.

Variațiile flexionare sunt schimbări ale morfemelor flexionare față de cele normate: *aeropoarte* (normat *aeroporturi*), *cenușei* (*cenușii*), *gazoane* (*gazonuri*), *girofare* (*girofaruri*).

Împrumuturi din engleză – ME: 19%

Împrumuturile din engleză ocupă, după cum se vede, aproape o cincime din totalul cuvintelor noi. Se pot distinge câteva domenii în care afluența de anglicisme este remarcabilă:

- nume de profesii: *advertiser, account manager, chief executive officer, chief head-hunter manager, sound-designer, web designer* ș.a.

- sport: *cross country, down hill, four-cross, full-contact, greencard, goal-keeper* ș.a.
- muzică: *backing vocal, boy-band, cover, dance* (adj.), *(deep) house, death-metal, downtempo, drum & bass, electric-gypsy, featuring, funk, heavy rotation* ș.a.
- modă: *baby doll, face-painting, fashion, fusion-fashion, gloss, hair-styling* ș.a.
- electronică/informatică: *dial-up, downloada, enter, gateway, hands-free, hi-tech* ș.a.
- alimentație: *fresh, junk-food, light* ș.a.

Există, firește, și împrumuturi din alte limbi, însă numărul lor este nesemnificativ. De asemenea sunt înregistrate și împrumuturi adaptate (adică împrumuturi care nu mai păstrează grafia din limba de origine), dar și ele sunt puțin numeroase. Numărul lor este cu atât mai mic cu cât s-a constatat o preferință a unor vorbitori de a folosi grafia originală (probabil ca pe un act de cultură) în ciuda faptului că cea adaptată a fost deja normată (a se vedea asemenea exemple la rubrica Variații ortografice). Exemple pot fi date din franceză: *biju*, (albastru, bleu) *ciel, connaisseur* ; chineză: *feng-shui, chi, thai-chi*; italiană: *dolce far niente, famiglia, silenzio stampa, Squadra Azzurra* ș.a.

Xenisme – Xenism: 0,5%

Deși semnificația generală a termenului *xenism* este aceea de împrumut neadaptat dintr-o limbă străină, noi deosebim xenismele de împrumuturile propriu-zise prin aceea că desemnează realități care nu sunt specifice zonei în care se vorbește limba-gazdă. Xenisme se găsesc, de pildă, în descrierile realităților din alte țări. asemenea exemple sunt *datsan* („curent budist”), *fatwa* („edict religios”), *kumara* („cultură agricolă din insula Paștelui”), *guerilleros* („luptători de gherilă”), *mzungus* („om din rasa albă”).

Altele – A: 5%

Pentru unele cuvinte nu se poate stabili etimologia, cum sunt cuvintele simple dialectale, argotice, cazuri ciudate, dificil de etichetat, dar care totuși sunt cuvinte nou folosite în presă. În această categorie am introdus termenii de specialitate din domeniul farmaceutic, medical etc. (*aflatoxină, botox, carmol*) sau cuvinte care nu au o etimologie evidentă (*amărăștean, amofarm, aszu, critifiction, flex* ș.a.).

5. Concluzii

Studiul de față atinge mai multe aspecte. În primul rând prezintă prima încercare de monitorizare sistematică a presei scrise românești grație resurselor electronice și informatice disponibile. Această cercetare poate constitui și o piatră de hotar pentru o monitorizare ulterioară precisă. Ne referim la faptul că obiectivul proiectului Neorom nu este atât cel de inventariere a neologismelor din presă, cât acela de inventariere a cuvintelor prezente în limbajul presei într-un anumit an precis. Or acest lucru nu se poate realiza, cât de cât exact, fără să existe un moment în care se trage linie și se spune „iată inventarul cuvintelor existente în limbă la acest moment”. În al doilea rând, acest studiu este important pentru lexicografia românească. S-a constatat că un număr destul de mare de cuvinte nu au fost încă înregistrate în dicționare în ciuda frecvenței mari de folosire și a vechimii lor. Multe dintre acestea reprezintă membri ai unor familii de cuvinte înregistrate parțial (de pildă apare în dicționar *minimalist*, dar nu și *minimalism*). De acum această deficiență poate fi corectată. Există și situații în care cuvintele din bază se găsesc în dicționare specializate pe domenii (argou, economic, informatică, sport

etc.) sau lingvistice (de sinonime, de antonime sau ortografice), însă faptul că aceste cuvinte apar în presă arată că ele sunt accesibile publicului larg și de aceea multe dintre ele ar trebui să se regăsească în dicționarele explicative, dacă depășesc un anumit prag de apariții. Trebuie spus că există și cuvinte care în mod surprinzător nu au fost înregistrate (precum *denim* sau *sclipici*). Chiar dacă nu orice cuvânt din baza de date Neorom se cuvine a fi introdus într-un dicționar al limbii române, studiul acesta este important și pentru lexicologie prin reflectarea creativității lingvistice, a tendințelor de evoluție a limbii precum și a raportului dintre limba română și celelalte limbi romanice implicate în proiect.

Mulțumiri. Doresc să mulțumesc d-nei Ioana Vintilă-Rădulescu, coordonatoarea părții române a proiectului Neorom, pentru sprijinul nemijlocit acordat în îmbunătățirea calității rezultatelor noastre. Mulțumesc de asemenea referenților acestei lucrări.

Referințe bibliografice

- Barbu, A.M. (2008). Romanian Lexical Data Bases: Inflected and Syllabic Forms Dictionaries. *Proceedings of Language Resources and Evaluation Conference - LREC 2008*, Marrakech și on-line <http://www.lrec-conf.org/proceedings/lrec2008/>.
- Barbu, A.M. (2010). Baza de date românească din cadrul platformei Neorom. *Limba română*, nr.2, București, sub tipar.
- Cabré, M. T. (coord.) (2004). *Metodología del trabajo en neología: criterios, materiales y procesos*. Observatori de Neologia, Papers de l'IULA. Sèrie Monografies, 9, Universitat Pompeu Fabra. Institut Universitari de Lingüística Aplicada, Barcelona.
- Coteanu, I., Forăscu, N., Bidu-Vrăncănu, A. (1985). *Limba română contemporană. Vocabularul*. Editura Didactică și Pedagogică, București.
- Stoichițoiu-Ichim, A. (2001). *Vocabularul limbii române actuale. Dinamică, influențe, creativitate*. Editura All, București.

Sigle

- DEX '98= *Dicționarul explicativ al limbii române*, Editura Univers Enciclopedic, București, 1998.
- DLRA = Zorela Creța, Lucreția Mareș, Zizi Ștefănescu-Goangă, Flora Șuteu, Valeriu Șuteu, *Dicționar al limbii române actuale* (ediția a II-a revăzută și adăugită), Editura Curtea Veche, București, 1998.
- DLRC = Vasile Breban, *Dicționarul limbii române contemporane*, Editura Științifică și Enciclopedică, București, 1980.
- DLRM= *Dicționarul limbii române moderne*, Editura Academiei, București, 1958.
- DN = Florin Marcu și Constant Maneca, *Dicționar de neologisme*, Editura Academiei, București, 1986.
- DOOM 1 = *Dicționar ortografic, ortoepic și morfologic al limbii române*, Editura Academiei Republicii Socialiste România, București, 1989.
- DOOM 2 = *Dicționar ortografic, ortoepic și morfologic al limbii române*, ediția a II-a, Editura Univers Enciclopedic, București, 2005.
- MDA = *Micul dicționar academic*, Editura Univers Enciclopedic, București, 2002.
- MDN = Florin Marcu, *Marele dicționar de neologisme*, Editura Saeculum, București, 2000.

EMOȚII ÎN CUVINTE: ELABORAREA RESURSEI MULTILINGVE

VICTORIA BOBICEV, VICTORIA MAXIM, TATIANA PRODAN,
NATALIA BURCIU, VICTORIA ANGHELUȘ

Universitatea Tehnică a Moldovei, Chișinău, Moldova

*vika@rol.md, maxivica@yahoo.com, tatiana.ursulenco@gmail.com,
natusicb@yahoo.com, lazu_vic@yahoo.com*

Rezumat

În lucrarea dată este descris procesul de creare a WordNet-Affect pentru limbile română și rusă. WordNet-Affect reprezintă o resursă lexicală elaborată în baza WordNet-ului englez ce conține informații despre emoțiile pe care le transmit cuvintele. WordNet-Affect este organizat în șase emoții de bază: *anger, disgust, fear, joy, sadness, surprise*.

Noi am tradus cuvintele din WordNet-Affect în limbile română și rusă și am creat o resursă lexicală aliniată engleză - română - rusă. Resursa dată este accesibilă gratuit în scopuri de cercetare.

1. Introducere

Creșterea spectaculoasă a tehnologiilor Web 2.0 permite fiecărui utilizator să participe activ la crearea de conținut web (bloguri, rețele sociale, chat-uri). Volumul de texte cu conținut emoțional bogat crește în progresie geometrică. Acest lucru face ca analiza subiectivă a textelor să fie de actualitate în mod special.

Până în prezent, analiza sentimentelor și studiile asupra afectelor cuvintelor s-au concentrat pe limba engleză. Un exemplu este sarcina SemEval-2007 de clasificare a textelor afective (Liu et al., 2003). De asemenea, cele mai multe resurse lexicale au fost create pentru limba engleză. De exemplu, SentiWordNet este o resursă lexicală pentru extragerea opiniilor care atribuie fiecărui synset al Wordnet-ului trei marcaje ale sentimentelor: pozitivitate, negativitate, obiectivitate (Esuli, Sebastiani, 2006).

Recent, cea mai mare creștere a utilizării Internetului a fost înregistrată pentru vorbitorii de limbi diferite de limba engleză: în anii 2000-2009, pentru limbile diferite de engleză, creșterea a depășit 690% în comparație cu 237% a creșterii pentru limba engleză¹.

În consecință, cantitatea informației textuale scrise în alte limbi decât engleza crește rapid ce sporește cererea pentru instrumentele de analiză automată a textelor și pentru resursele lexicale pentru alte limbi, diferite de engleză. Elaborarea instrumentelor a progresat pentru limbile vest-europene (franceză, germană) și asiatice (japoneză, chineză, arabă) (Edmonds, 2004), se simte nevoia de astfel de resurse și pentru limbile est-europene. Pentru a umple acest gol, am elaborat o resursă lexicală, în baza WordNet-Affect, prin traducerea în română și rusă, redactarea synset-urilor traduse și alinierea acestora la sursa engleză.

¹ <http://www.internetworldstats.com/stats7.htm>

2. WordNet-Affect

WordNet-Affect² este o resursă lexicală ce conține informații despre emoțiile pe care le transmit cuvintele. În comparație cu WordNet-ul complet, WordNet-Affect este o resursă lexicală mică, dar valoroasă pentru adnotarea sa afectivă.

a#01943022 awed awestruck awestricken in_awe_of

Figura 1: Un synset din WordNet-Affect.

WordNet-Affect (Strapparava, Valitutti, 2004) a fost creat pornind de la WordNet DOMAINS (Magnini, Cavaglia, 2002). WordNet-Affect produce o ierarhie suplimentară a etichetelor domeniilor afective, independent de ierarhia domeniilor, cu care sunt adnotate synset-urile ce reprezintă concepte afective. „Cuvintele afective” sunt considerate a fi cuvintele ce au „conotație emoțională” (Ortony, 1987). Există cuvinte ce nu descriu direct unele emoții (de exemplu, bucurie, tristețe sau frică), dar, de asemenea, sunt legate de emoții precum cuvintele ce descriu stările mintale, stările fizice sau trupești, trăsăturile de personalitate, comportamente, atitudini și sentimente (cum ar fi plăcere sau durere).

Tabel 1: Seturile de date ale cuvintelor afective.

Clasele	#synset-uri	%synset-uri	#cuvinte	%cuvinte
anger	128	21.0	318	20.7
disgust	20	3.3	72	4.7
fear	83	13.5	208	13.5
joy	228	37.2	539	35.1
sadness	124	20.3	309	20.1
surprise	29	4.7	90	5.9
Total	612	100.0	1536	100.0

Colecția de synset-uri din WordNet-Affect utilizată în lucrarea dată a fost furnizată ca o resursă pentru SemEval-2007 „Affective Text”³. Această sarcină a fost axată pe adnotarea textelor prin etichete afective (Strapparava, Mihalcea, 2008). Toate synset-urile din WordNet-Affect sunt adnotate utilizând șase etichete de categorii emoționale: *joy*, *fear*, *anger*, *sadness*, *disgust*, *surprise* (Strapparava et al., 2006). Această alegere a celor șase emoții rezultă din cercetarea psihologică a emoțiilor umane exprimate non-verbal (Ekman, 1992).

Datele despre resursa inițială sunt descrise în tabelul 1. Toată resursa este prezentată în șase fișiere denumite conform celor șase emoții. Fiecare fișier conține o listă de synset-uri; un synset pe un singur rând. Un exemplu de synset este prezentat în figura 1.

Prima literă din rând indică partea de vorbire; aceasta este urmată de numărul synset-ului și apoi sunt enumerate toate cuvintele din synset. Această reprezentare a fost simplă și ușoară pentru prelucrarea ulterioară. În synset-uri inițiale este un număr mare de combinații de cuvinte, sintagme și expresii. Un exemplu poate fi văzut în figura 1: „in_awe_of”. Aceste părți ale synset-urilor au generat probleme în timpul traducerii.

² În scopuri de cercetare, WordNet-Affect este accesibil la cerere la <http://wndomains.itc.it>

³ <http://www.cse.unl.edu/~rada/affectivetext>

3. Elaborarea WordNet-ului pentru limbile română și rusă

WordNet-ul Român a fost creat de către Universitatea Alexandru Ioan Cuza din Iași în cadrul proiectului european Balkanet (Tufiş et al., 2006). După ce proiectul Balkanet s-a încheiat, Institutul de Cercetări pentru Inteligență Artificială, de la Academia Română a continuat să lucreze cu WordNet-ul Român și, în prezent, acesta conține 33151 synset-uri de substantive, 8929 synset-uri de verbe, 851 synset-uri de adjective și 834 synset-uri de adverbe (Tufiş et al., 2008). Acesta poate fi accesat prin intermediul interfeței online a MultiWordNet⁴ în care WordNet-urile pentru câteva limbi sunt aliniate la WordNet-ul Princeton.

În lucrarea dată, mai întâi de toate am verificat synset-urile din WordNet-Affect folosind interfața online a MultiWordNet. Astfel, am copiat la setul nostru toate synset-urile care deja există în WordNet-ul Român și nu am procesat aceste synset-uri suplimentar. Ca rezultat, în WordNet-ul Român au fost găsite 166 synset-uri, majoritatea dintre ele fiind pentru substantive și verbe. Adjectivele și adverbele sunt reprezentate mai puțin. Statistica privind synset-urile deja existente în română este prezentată în tabelul 2.

Tabel 2: Seturile de date a synset-urilor deja existente în Wordnet-ul Român.

Clasele	# synset-uri în WordNet-Affect	# synset-uri din WordNet Român	% synset-uri din WordNet Român
anger	116	35	30.1
disgust	17	7	41.1
fear	76	25	32.8
joy	209	63	30.0
sadness	98	24	24.7
surprise	26	12	46.1
Total	542	166	30.6

Pentru limba rusă situația este complet diferită. Au fost făcute mai multe încercări pentru a crea WordNet-ul pentru limba rusă. RussNet este un proiect pentru crearea tezaurului rusesc electronic (Azarova et al., 2002). Un proiect alternativ al versiunii rusești a WordNet-ului este WordNet-ul Rus (Balkova et al., 2004). Ambele proiecte sunt necomerciale. De asemenea, există două proiecte comerciale menite să elaboreze WordNet-uri pentru limba rusă: RuThes este un tezaur informațional utilizat în UIS RUSSIA⁵ și proiectul privind WordNet-ul Rus creat de către grupul companiei Novosoft⁶. Din păcate, puțină informație și încă mai puține resurse sunt accesibile în mod gratuit.

4. Crearea WordNet–Affect pentru limbile română și rusă

Pentru crearea acestor două seturi de date, am parcurs trei etape: (1) traducere automată; (2) înlăturarea traducerilor irelevante; (3) generarea synset-urilor pentru română și rusă.

⁴ <http://multiwordnet.itc.it/english/home.php>

⁵ <http://www.cir.ru>

⁶ <http://research-and-development.novosoft-us.com>

4.1 Traducerea automată

Traducerea a fost făcută automat utilizând dicționare bilingve. Am folosit dicționarul electronic Român–Englez ROMEN de la PRIMASOFT⁷. Dicționarul constă din următoarele compartimente: Englez–Român, Român–Englez, Englez–Rus, Rus–Englez, fiecare conținând mai mult de 200 000 de intrări. Noi am utilizat doar compartimentele în care limba engleză este limba sursă. În dicționare au fost combinații de cuvinte, sintagme și expresii pe care le-am folosit în limbile țintă. Pentru traducerea automată, dicționarul a fost organizat într-o listă de cuvinte sursă urmate de traducerile țintă. Un exemplu de intrare a dicționarului este prezentat în figura 2.

Joy
Dicționar general:
noun: bucurie; confort; fericire; plăcere; tihnă; veselie;
voioșie;
verb: a bucura; a înveseli;

Figura 2: Un exemplu de intrare a dicționarului.

La această etapă, scopul nostru a fost să obținem cât mai multe cuvinte afective posibile pentru analiză. Pentru aceasta, am tradus fiecare cuvânt din synset-urile WordNet-Affect. Am decis să excludem din synset-urile engleze toate combinațiile de cuvinte, sintagmele și expresiile, deoarece ele nu pot fi traduse automat. Figura 3 prezintă un exemplu de synset tradus, obținut după această etapă. După cum se vede în exemplu, pentru traducerea română, de asemenea am obținut combinații de cuvinte care erau în dicționar: „a avea gust”, „a degusta (un aliment)”.

Unele elemente ale synset-urilor nu au fost traduse. Acestea pot fi divizate în patru grupe. (1) Combinații de cuvinte, sintagme și expresii pe care le-am înlăturat intenționat din synset-urile engleze înainte de traducere. (2) Variații în ortografia aceluiași cuvânt; de exemplu, „jubilance”, „jubilancy” – primul cuvânt a fost tradus, al doilea nu a fost găsit în dicționar. (3) Cuvinte care au fost formate cu ajutorul sufixelor „ness”, „less”, „ful” (de exemplu „heartlessness”); este puțin probabil ca acestea să apară în dicționare, la fel ca și adverbele formate cu sufixul „ly”. (4) Cuvinte care nu au fost traduse din cauza limitărilor dicționarului utilizat. În timp ce WordNet poate fi, în mod rezonabil, menționat ca unul din cele mai mari dicționare engleze, dicționarul nostru bilingv este destul de modest. Tabelul 3 arată procentajul cuvintelor care nu au fost traduse. Media în procente a cuvintelor netraduse a fost 21%.

⁷ http://www.primasoft.biz/romen_eng.php

EMOȚII ÎN CUVINTE: ELABORAREA RESURSEI MULTILINGVE

Tabel 3: Numărul și procentajul cuvintelor netraduse.

Clase	# de cuvinte engleze	# de cuvinte traduse	# de cuvinte netraduse	% de cuvinte netraduse
anger	318	248	70	22.0
disgust	72	60	13	18.0
fear	208	162	47	22.5
joy	539	420	119	22.0
sadness	309	246	63	20,5
surprise	90	72	18	21.0
Total	1536	1208	330	21.0

A doua grupă de cuvinte nu prezintă o problemă, dar prima, a treia și a patra au trebuit traduse manual. Aceasta s-a făcut la etapa a treia.

05573914 n:

preference =

preferință

penchant =

înclinație

slăbiciune

predilection =

predilecție

taste =

a avea gust

a gusta, a cunoaște

a gusta; a degusta (un aliment)

degustare

fărămă, bucățică, îmbucătură (de)

gust

înclinație, preferință

Figura 3: Un exemplu de traducere a unui synset.

4.2 Înlăturarea traducerilor irelevante

Multe cuvinte din synset-urile engleze au câteva sensuri. Este evident că traducerea automată a generat toate traducerile posibile pentru toate sensurile. Noi am ales doar o singură traducere care era relevantă pentru semnificația synset-ului. Traducerea relevantă a fost selectată manual. Am înlăturat toate traducerile a căror sensuri nu aveau legătură cu emoția. De exemplu, cuvântul „taste” în synset-ul cu înțelesul de „preference” are câteva semnificații, însă doar ultima din lista traducerilor posibile avea legătură cu sensul comun al synset-ului. Exemplul este prezentat în figura 3. Astfel, noi am înlăturat toate traducerile cu excepția ultimei.

Deoarece am tradus fiecare cuvânt în mod separat, am obținut o mulțime de duplicate care la fel au trebuit eliminate. Am tras atenția, de asemenea, la corespondența părților de vorbire. În multe cazuri, a fost destul de dificil, în special pentru substantivele deja menționate, formate cu ajutorul sufixelor, de exemplu, „plaintiveness” sau „uncheerfulness”.

4.3 Generarea synset-urilor române și ruse

Toate cuvintele dintr-un synset reprezintă un concept, un singur sens. Scopul etapei a treia a fost de a găsi traducerea adecvată, exact corespunzătoare acestui sens. La această etapă, a trebuit mai întâi să atașăm explicația în engleză a fiecărui synset. Aceasta a determinat claritatea în semnificația synset-ului pentru traducători. După ce au fost adăugate explicațiile la synset-uri, întregul set a fost dat la trei traducători care au lucrat independent. Sarcina lor a fost dublă: (1) să elimine traducerile care, din punctul lor de vedere, nu erau relevante semnificației synset-ului descrise de explicația acestuia; (2) să adauge cât mai multe sinonime posibile relevante la synset-urile române și ruse. Astfel, sarcina lor a fost să verifice echivalența sensurilor synset-urilor engleze, ruse și române. De asemenea, au trebuit să traducă cuvintele care au rămas netraduse la prima etapă. Pentru traducere au fost folosite dicționare online.

Dicționarele bilingve române utilizate:

- <http://hallo.ro>,
- <http://dictionar.netflash.ro>,
- <http://www.ectaco.co.uk/English-Romanian-Dictionary>;

Dicționar explicativ român: <http://dexonline.ro/>.

Dicționarele bilingve ruse utilizate:

- <http://en.bab.la>,
- <http://dictionary.babylon.com>,
- <http://russianlessons.net/dictionary/dictionary.php>;

Dicționare explicative ruse:

- <http://slovo.freecopy.ru/>,
- <http://slovari.yandex.ru/dict/ushakov>.

Această etapă a fost cea mai dificilă care a necesitat cel mai mult lucru. Multe synset-uri engleze aveau aproape aceleași sensuri, doar cu careva nuanțe. În unele cazuri, synset-urile conțineau cuvinte arhaice, care nu erau găsite în dicționare. După cum s-a menționat mai sus, am încercat să evităm combinațiile de cuvinte, sintagmele și expresiile. Cu toate acestea, în unele cazuri, sensul exact al unui synset englez putea fi reprezentat doar de o singură combinație de cuvinte române sau ruse. În unele cazuri, chiar și synset-ul englez era prezentat de combinații de cuvinte. De exemplu, n#05591681 stage_fright. Un alt exemplu conține un cuvânt german: n#05600844 world-weariness Weltschmerz. În așa cazuri, nu am obținut traducerea corespunzătoare. În unele cazuri, câteva synset-uri engleze au fost traduse în aceleași cuvinte române sau ruse deoarece nu am putut reflecta în limba țintă nuanțele sensurilor limbii sursă.

Referindu-se la problema cu sufixe, de exemplu, cuvintele „weepiness”, „plaintiveness”, „mournfulness”, „ruthfulness”, cu greu pot fi găsite în dicționare. Pentru a rezolva această problemă, am eliminat sufixele și am căutat rădăcinile cuvintelor menționate în dicționare disponibile. În acest fel, am putut găsi sensul cuvintelor și, prin adăugarea afixelor necesare, au fost create echivalentele în

română și rusă. De exemplu, pentru a găsi traducerea adecvată pentru cuvântul „mournfulness”, am căutat în dicționar cuvântul „mournful”. Rezultatul pentru limba română este „îndoliat” și pentru rusă „траурный”. Având în vedere că cuvântul „mournfulness” este un substantiv, am transformat adjectivele obținute în substantive. Astfel, echivalentul românesc este „doliu” și în rusă - „тпайр”.

Cu toate acestea, cele mai multe probleme au apărut la alinierea adjectivelor. De exemplu, pentru eticheta emoțională „sadness”, multe synset-uri adjectivale traduse în limba rusă conțin cuvintele „грустный” și „печальный”. Prin urmare, pentru synset-uri adjectivale diferite obținem traduceri destul de asemănătoare.

4.4 Acordul între traducători

Având ca rezultatul traducerii seturi de sinonime, nu am putut folosi măsurile standard pentru acordul între traducători. Astfel, acordul a fost calculat după cum urmează. Dacă A este o mulțime de cuvinte selectate de către primul traducător pentru un synset și B este o mulțime de cuvinte selectate de al doilea traducător pentru același synset, acordul între traducători $IntAgr$ este egal cu raportul dintre numărul de cuvinte în intersecția mulțimilor A și B și numărul de cuvinte în uniunea mulțimilor A și B:

$$IntAgr = (A \cap B) / (A \cup B) \quad (1)$$

De exemplu, dacă un traducător a format un synset din trei cuvinte w_k, w_l și w_m , iar al doilea traducător a format același synset din patru cuvinte w_k, w_l, w_m și w_n , și primele trei cuvinte sunt aceleași, atunci $A = (w_k \ w_l \ w_m)$, $B = (w_k \ w_l \ w_m \ w_n)$, $A \cap B = (w_k \ w_l \ w_m)$, $A \cup B = (w_k \ w_l \ w_m \ w_n)$, numărul de cuvinte din intersecția A și B ar fi egal cu 3, numărul de cuvinte din uniunea A și B ar fi egal cu 4 și, prin urmare, acordul între traducători - ar fi $3 / 4 = 0.75$.

Spre exemplu synset-ul „a#01195320 friendly” a fost tradus de primul traducător ca „prietenos prietenesc amical”, de al doilea traducător ca „binevoitor amical prietenos”, și de al treilea, ca „binevoitor prietenesc prietenos”. Pentru primul și al doilea traducător intersecția traducerilor era de două cuvinte: „prietenos amical” și uniunea traducerii a fost patru cuvinte „prietenos prietenesc binevoitor amical”. Acordul între traducători, în acest caz a fost de $2 / 4 = 0.5$. Pentru al doilea și al treilea traducător, intersecția traducerilor era de două cuvinte: „prietenos binevoitor”, și uniunea - patru cuvinte „prietenos prietenesc binevoitor amical”. Prin urmare, acordul este același: 0.5. Pentru primul și al treilea traducător acordul, de asemenea, este același: 0.5. Toți cei trei traducători au în comun doar un singur cuvânt „prietenos” și uniunea de traduceri a constat din patru cuvinte. Astfel, acordul a fost de $1 / 4 = 0.25$.

Tabelul 4 prezintă valorile medii ale acordului între traducători. Cei trei traducători sunt prezentați ca T1, T2 și T3.

Tabel 4: Acordul între traducători.

Perechile de traducători	Acordul între traducători
Date pentru limba rusă	
T1 – T2	0.57
T2 – T3	0.61
T1 – T3	0.59
Toți	0.29
Date pentru limba română	
T1 – T2	0.58
T2 – T3	0.57
T1 – T3	0.67
Toți	0.32

Unele synset-uri aveau acordul egal cu unu, spre exemplu, în synset-ul „a#00863650 euphoriant”, toți trei traducători au tradus cuvânt respectiv ca „euforizant”. Cu toate acestea, pentru majoritatea synset-urilor, traducătorii au furnizat mai multe traduceri diferite, dar nu multe dintre aceste traduceri au fost comune pentru toți traducătorii. În unele synset-uri traduse, nu a fost găsit nici un cuvânt comun în traduceriile celor trei traducători. De exemplu, pentru synset-ul „a#00670851 gladdened exhilarated”, cele trei traduceri au fost „bucurat înveselit înviorat bine_dispus”, „bucuros vesel voios încântat bine_dispus” și „bucurat voios bucuos înveselit”. Nu a fost găsit nici un cuvânt comun pentru toate cele trei traduceri. Astfel, ne-am decis să formăm synset-uri din cuvintele care au apărut în cel puțin două dintre cele trei variante de traduceri. În acest fel, am format synset-uri finale. De exemplu, synset-ul „a#01195320 friendly” a fost tradus ca „prietenos prietenesc binevoitor amical”, pentru că toate aceste cuvinte au apărut de cel puțin două ori în traduceri. Synset-ul „a#00670851 gladdened exhilarated” a fost tradus ca „bucurat înveselit bine_dispus bucuos voios”.

Tabelul 5 conține date privind numărul final de cuvinte în traduceri pentru fiecare din cele șase emoții din WordNet-Affect .

Tabel 5: Seturile de date a cuvintelor afective din limba română și rusă.

Clasele	#synset-urile	# Cuvinte ruse	# Cuvinte române
anger	116	393	330
disgust	17	73	60
fear	76	327	248
joy	209	765	641
sadness	98	437	364
surprise	26	129	87
Total	542	2199	1869

Trebuie de menționat că în sursa WordNet-Affect englez au fost găsite unele synset-uri duplicate. Am eliminat toate aceste repetiții și numărul de synset-uri în sursa noastră este mai mic. În plus, au fost depistate unele diferențe mici dintre WordNet-Affect,

MultiWordNet și versiunea online a WordNet pentru că WordNet-Affect s-a creat în baza WordNet versiunea 1.6, MultiWordNet utilizează versiunea 2.0 a WordNet-ului și versiunea online a WordNet-ului este 3.0. În ciuda numărului mai mic de synset-uri, numărul de cuvinte în seturile pentru limbile română și rusă este mai mare decât în limba engleză. Acest lucru se datorează tendinței noastre de a colecta în resursele noastre cât mai multe cuvinte posibile. Scopul nostru este de a folosi această resursă pentru metodele statistice de recunoaștere a emoțiilor în text.

5. Concluzii și planuri de viitor

Acest articol descrie procesul de creare a resursei lexicale multilingve WordNet-Affect ce conține partea engleză, partea rusă și partea română aliniată la nivel de synset. WordNet-Affect este o resursă lexicală creată în baza WordNet-ului Princeton care codifică informațiile despre emoțiile pe care le transmit cuvintele. Acesta este organizat în șase emoții de bază: *anger, disgust, fear, joy, sadness, surprise*. WordNet-Affect este o resursă lexicală comparativ mică, dar valoroasă pentru adnotările sale afective.

În lucrare este descris procesul traducerii synset-urilor WordNet-Affect în limba română și rusă și creării WordNet-Affect aliniat Engleză - Română - Rusă. Resursa poate fi folosită pentru recunoașterea automată a emoțiilor și afectelor în text. Aceasta poate fi obținută gratuit în scopuri de cercetare, pe <http://lilu.fcim.utm.md>. Resursa este încă în dezvoltare. Prima versiune bazată pe WordNet-Affect a fost lansată în august 2009; a doua, lansată în octombrie 2009, este aliniată cu WordNet-ul român. Resursa a fost deja folosită în (Sokolova, Bobicev, 2009) și aceasta este numai una dintre multiplele posibilități de utilizare a seturilor de cuvinte.

În viitor plănuim să mărim resursa noastră. Recent am obținut WordNet Domains care este distribuit cu WordNet-Affect. În resursa dată sunt incluse două versiuni ale WordNet-Affect original englez. Prima versiune a fost obținută semiautomat și conține aproximativ 1 500 de synseturi. A doua versiune a fost creată după verificarea manuală a primei versiuni și excluderea synset-urilor care aveau o legătură mai slabă cu emoția respectivă. Versiunea aceasta conține 914 synset-uri. Intenția noastră este de a traduce toate 914 synset-uri și de a alinia acestea la WordNet român. În afară de aceasta, vom crea resurse 'bag-of-words' pentru utilizarea imediată în sisteme de recunoaștere a emoțiilor și afectelor în text.

Referințe bibliografice

- Azarova, I., Mitrofanova, O., Sinopalnikova, A., Yavorskaya, M., Oparin I. (2002). Russnet: Building a lexical database for the russian language. *Workshop on Wordnet Structures and Standardization and How this affect Wordnet Applications and Evaluation*, Las Palmas, pp. 60-64.
- Balkova V., Suhonogov A., Yablonsky S.A. (2004). Russian WordNet. From UML-notation to Internet/Intranet Database Implementation. *Second International WordNet Conference, GWC 2004*, Brno, Czech Republic, 31-38.
- Edmonds, P. (2002). Introduction to Senseval. *ELRA Newsletters*, 7(3), 337-344.

- Ekman, P. (1992). An argument for basic emotions. *Cognition and Emotion*, vol. 6(3-4), 169–200.
- Liu, H., Lieberman, H., Selker, T. (2003). A model of textual affect sensing using real-world knowledge. *ACM Conference on Intelligent User Interfaces*, Miami, Florida, USA, pp. 125-132.
- Strapparava, C., Mihalcea, R. (2008). Learning to identify emotions in text. *ACM Symposium on Applied Computing*, Fortaleza, Brazil, 556-1560.
- Strapparava, C., Valitutti, A. (2004). Wordnet-affect: an affective extension of wordnet. *4th International Conference on Language Resources and Evaluation*, 1083–1086.
- Strapparava, C., Valitutti, A., Stock, O. (2006). The affective weight of the lexicon. *5th International Conference on Language Resources and Evaluation (LREC 2006)*, Genoa, Italy, 474–481.
- Tufiş, D., Mititelu, B., Bozianu, L., Mihaila, C. (2006). Romanian wordnet: New developments and applications. *3rd Conference of the Global WordNet Association*, Korea, 337–344.
- Tufiş, D., Ion, R., Bozianu, L., Ceaşu, A., Ştefănescu, D. (2008). Romanian Wordnet: Current State, New Applications and Prospects. *4th Global WordNet Conference, GWC-2008*, University of Szeged, Hungary, 441-452.
- Esuli, A., Sebastiani, F. (2006). SentiWordNet: A Publicly Available Lexical Resource for Opinion Mining. *5th International Conference on Language Resources and Evaluation (LREC 2006)*, Genoa, Italy, 417-422.
- Magnini, B., Cavaglia G. (2002). Integrating subject field codes into Wordnet. *Second International Conference on Language Resources and Evaluation (LREC 2002)*, Athens, Greece, 1413—1418.
- Ortony, A., Clore, G. L., Foss, M. A. (1987). The psychological foundations of the affective lexicon. *Journal of Personality and Social Psychology*, American Psychological Association, vol. 53, 751–766.
- Sokolova M., Bobicev V. (2009). Classification of Emotion Words in Russian and Romanian Languages. *RANLP-2009 conference*, Borovets, Bulgaria, 415-419.

CAPITOLUL 3

APLICAȚII ALE TEHNOLOGIILOR LINGVISTICE TEXTUALE

SISTEM ÎNTREBARE-RĂSPUNS ANTRENABIL PENTRU LIMBA ROMÂNĂ

DAN ȘTEFĂNESCU, RADU ION, ALEXANDRU CEAUȘU, DAN TUFIȘ, ELENA IRIMIA, VERGINICA BARBU-MITITELU

Institutul de Cercetări pentru Inteligență Artificială, Academia Română

{danstef, radu, aceausu, tufis, elena, vergi}@racai.ro

Rezumat

Lucrarea prezintă un sistem întrebare-răspuns dezvoltat la *Institutul de Cercetări pentru Inteligență Artificială* – ICIA în cadrul unui proiect național și evaluat independent în contextul competiției europene CLEF. Evaluarea a fost realizată în cadrul exercițiului *ResPubliQA* pentru limba română. Este descris modul de combinare a diferiților factori de relevanță pe baza cărora sistemul identifică paragrafele cele mai relevante ca răspunsuri la întrebările formulate în limbaj natural. Sistemul este disponibil on-line pe pagina de servicii web a ICIA. El este însă complet antrenabil, funcționalitatea sa fiind independentă de registrul lingvistic ce caracterizează datele de antrenare.

1. Introducere

Cercetările privind prelucrarea automată a limbajului natural (PLN), domeniu central al inteligenței artificiale, produc rezultate cu impact din ce în ce mai mare în societatea globalizată de fenomenul Internet. Sistemele de întrebare-răspuns (ÎR) în limbaj natural, în vogă în perioada anilor '70-80, au revenit în centrul cercetărilor PLN dar, de data aceasta având ca obiectiv identificarea răspunsurilor la întrebări arbitrare în spații de căutare incomparabil mai mari, la limită întregul web. Conținutul informațional al acestui spațiu virtual este atât de mare încât se consideră că orice solicitare rațională de informație își poate găsi măcar un răspuns pe Internet. Evaluarea calității răspunsurilor și respectiv asigurarea găsirii lor sunt însă probleme de cercetare, pentru care abordările tradiționale au devenit insuficiente. Astfel, implementările motoarelor de căutare moderne recurg din ce în ce mai mult la tehnici PLN, acestea fiind utilizate în toate etapele fluxului de prelucrare, începând de la nivelul specificării întrebării și până la extragerea fragmentului de text relevant dintr-unul sau mai multe documente. Odată cu utilizarea tehnicilor PLN au apărut și campaniile de evaluare în domeniul regăsirii inteligente a informației. Acestea constituie astăzi priorități ale cercetării de avangardă dedicată spațiului digital al cunoașterii. Ele au fost organizate mai întâi în SUA (*MUC-Message Understanding Conference*, *TREC-Text Retrieval Conference*, *DUC-Document Understanding Conference*, devenită *TAC-Text Analysis Conference*). În Europa, manifestarea similară este CLEF (Cross Language Evaluation Forum), ajunsă în anul 2009 la a 10-a ediție. Având ca subiect al analizei în primul rând limbile Uniunii Europene, începând cu anul 2006 limbile europene „cu resurse electronice limitate” (română, bulgară, cehă, greacă, portugheză etc.) au devenit „subiecte” de concurs.

În acest context, proiectul național SIR-RESDEC (lansat în 2007) a răspuns unei priorități europene, propunându-și realizarea unui sistem de ÎR în limbaj natural la nivelul celor mai avansate sisteme ale cercetării internaționale. Consorțiul SIR-RESDEC format din cercetători de la ICIA, UAIC și ICI și-a concentrat eforturile în

direcția realizării unor sisteme de ÎR pentru limbile română și engleză, în domeniul legislației Uniunii Europene și respectiv în domeniul geneticii umane. Grupul de cercetare al ICIA și-a concentrat eforturile în direcția realizării sistemului de ÎR pentru limba română în domeniul legislativ, având la dispoziție corpusul JRC-Acquis (Steinberger et al., 2006). În restul articolului de față vom descrie acest sistem.

2. Considerente preliminare

Pentru a testa performanțele sistemului, ICIA s-a înscris în anul 2009 în competiția QA@CLEF la secțiunea *ResPubliQA*¹ (Peñas et al., 2009), urmând tradiția participărilor la competițiile CLEF încă din 2006. Sarcina sistemelor dezvoltate de echipele înscrise la *ResPubliQA* a fost să identifice automat paragrafe relevante pentru răspunsuri la întrebări formulate în limbaj natural, în domeniul juridic acoperit de corpusul paralel de lucru al competiției. Pentru prima oară, evaluările sistemelor de întrebare-răspuns în limbaj natural au putut fi comparate interlingual, întrucât întrebările de test (500) au fost aceleași în 8 limbi (bască, bulgară, engleză, franceză, germană, italiană, română și spaniolă) răspunsurile trebuind a fi căutate în corpusul paralel (JRC-Acquis) al legislației europene „Acquis Communautaire” disponibil în toate limbile Uniunii Europene. Alinierea la nivel de paragraf a corpusului pentru toate limbile implicate a oferit posibilitatea evaluării răspunsurilor sistemelor indiferent de limba de interogare. În plus, tot în premieră, organizatorii ResPubliQA au calculat, pentru fiecare limbă, performanțele unui sistem de regăsire documentară (RD) de ultimă generație, dar fără componente de PLN. S-a urmărit în acest mod evaluarea cantitativă a rolului tehnologiilor de prelucrare a limbajului natural față de tehnicile standard utilizate în regăsirea informațiilor. Notând cu A_{IR} acuratețea sistemului de ÎR (v. Secțiunea 6) și cu A_{RD} acuratețea sistemului de RD, atunci raportul $M = \frac{A_{\text{IR}}}{A_{\text{RD}}}$ cuantifică meritul tehnicilor de prelucrare a limbajului natural. O cifră de merit supraunitară semnifică faptul că prelucrarea limbajului natural îmbunătățește performanța unui sistem de regăsire documentară. Deși pare intuitiv ca pentru orice sistem de ÎR cifra de merit M să fie supraunitară, organizatorii ResPubliQA au constatat că dintre 28 de sisteme evaluate doar 14 au avut o cifră de merit supraunitară (v. Table 10 în (Peñas et al., 2009)).

Cu experiența dobândită în competițiile CLEF precedente (Pușcașu et al., 2007; Tufiș et al., 2008c; Ion et al., 2009a) și cerințele specifice ale competiției din 2009, obiectivele tehnice principale au fost perfecționarea modului de regăsire a paragrafelor relevante și respectiv a celui de reordonare a acestora, pe baza unei analize complexe a relevanței paragrafelor candidat. Un impact semnificativ l-a constituit implementarea unei metode similare metodei de optimizare MERT (Och, 2003), în cadrul etapei de reordonare a paragrafelor. Am păstrat abordarea din anii precedenți în ceea ce privește construcția sistemului, însă diferitele module care îl alcătuiesc au fost implementate ca servicii și/sau aplicații Web (Tufiș et al., 2008b):

- *serviciul de analiză a întrebării*² cu ajutorul căruia se clasifică fiecare întrebare atașându-i-se o etichetă ce indică tipul de răspuns pe care acea întrebare îl cere;

¹ <http://celct.isti.cnr.it/ResPubliQA/>

² <http://shadow.racai.ro/JRCACQCWebService/Service.aspx?WSDL>

- *serviciile de generare a cererilor*³ cu ajutorul căruia o întrebare în limbaj natural este transformată în interogări în limbaj formal compatibile cu motorul de căutare;
- *serviciul de interogare a motorului de căutare*⁴ se ocupă de partea de regăsire a paragrafelor relevante pentru o interogare în limbaj formal furnizată la intrare;
- *modulul de reordonare a paragrafelor* preia rezultatele furnizate de motorul de căutare sub forma unei liste de paragrafe, calculează scoruri suplimentare de relevanță pentru fiecare paragraf și, în funcție de o interpolare liniară (obținută aplicând o optimizare de tip MERT) a acestora, asignează scoruri paragrafelor. Paragraful sau paragrafele cu scorurile cele mai ridicate sunt întoarse utilizatorilor.

În faza de indexare a corpusului JRC-Acquis, am luat în considerare doar textul propriu-zis al documentelor, acesta fiind în prealabil preprocesat cu ajutorul TTL (Ion, 2007) dezvoltat la ICIA. Textul a fost segmentat la nivel de unitate lexicală în funcție de terminologia Eurovoc, adnotat la parte de vorbire, și lematizat.

3. Identificarea terminologiei

Având în vedere caracterul juridic specializat al corpusului de lucru, o etapă importantă a fost identificarea și tratarea ca unități lexicale a unor anumite expresii sau termeni multi-cuvânt. Acest lucru a fost realizat cu ajutorul tezaurului multilingv Eurovoc⁵ (Ștefănescu and Tufiș, 2006). Descriptorii Eurovoc sunt termeni tehnici care trebuie să apară cu consecvență în toate documentele juridice pentru toate limbile implicate. Recunoașterea acestora atât în faza de analiză a întrebărilor cât și faza de preprocesare a documentelor (ce trebuie ulterior indexate) devine esențială pentru performanțele oricărui sistem QA pe acest corpus.

În consecință am realizat un modul care, după preprocesarea corpusului, recunoaște termenii Eurovoc și generează unități lexicale corespunzătoare (în acest fel, unui termen multi-cuvânt îi corespunde o singură unitate lexicală). Etapa de identificare a terminologiei se desfășoară de-a lungul a 6 etape: (i) termenii din Eurovoc sunt identificați în corpus, în forma lor de tezaur; (ii) pentru fiecare formă ocurență identificată la pasul 1 se extrage secvența de leme implicată și se adaugă unui inventar; (iii) se identifică în corpus toate ocurențele secvențelor de leme din inventarul construit la pasul 2; (iv) termenii astfel identificați sunt aduși la forma ocurență din corpus; (v) fiecărui termen identificat i se asignează un descriptor morfo-sintactic – cum termenii sunt de fapt grupuri nominale, descriptorul asignat termenului este același cu descriptorul centrului grupului nominal; (vi) fiecărei ocurențe a unui termen i se asignează ca leme descriptorul corespunzător Eurovoc. De exemplu, termenul *adunare parlamentară* apare în corpus cu formele flexionate: *adunarea parlamentară*, *adunările parlamentare*, *adunărilor parlamentare*. Toate aceste unități lexicale primesc ca leme descriptorul Eurovoc *adunare parlamentară*, iar ca descriptor morfosintactic, descriptorul corespunzător centrului grupului (i.e., *adunare*, *adunarea*, *adunările*, etc.).

³ <http://shadow.racai.ro/QADWebService/Service.aspx?WSDL>

⁴ <http://www.racai.ro/webservices/search.aspx?WSDL>

⁵ <http://en.wikipedia.org/wiki/Eurovoc>

4. Clasificarea automată a paragrafelor și întrebărilor

Specificațiile competiției au definit 5 tipuri de întrebări posibile: (i) *factoid* (factual) – întrebări care cer ca răspuns persoane, locații, instituții, momente în timp, etc.; (ii) *definition* (definiție) – întrebări care cer ca răspuns o definiție; (iii) *procedure* (procedură) – întrebări care cer ca răspuns o procedură juridică; (iv) *reason* (motiv) – întrebări care cer ca răspuns un motiv, o cauză; (v) *purpose* (scop) – întrebări care cer ca răspuns un scop, un obiectiv. Numărul redus de clase și faptul că răspunsul corect pentru o întrebare nu se poate întinde pe mai multe paragrafe au condus la ideea clasificării paragrafelor în funcție de probabilitatea lor de a răspunde la un tip de întrebare sau altul. Etichetarea unui paragraf cu tipul de întrebare la care acel paragraf ar putea răspunde cel mai bine, oferă, în mod evident, posibilitatea reducerii complexității etapei de identificare a documentelor relevante, pe care o vom numi etapă de *regăsire documentară*. Acolo unde tipul întrebării este corect identificat, răspunsul este căutat, în principal, în paragrafele de tip identic cu cel al întrebării. Acest lucru a presupus construcția unui modul de clasificare a întrebărilor care a fost optimizat și antrenat pentru tipurile de întrebări asociate corpului de tip juridic așa cum vom arata mai jos.

Problema de clasificare a paragrafelor este similară cu cea a selecției propozițiilor dintr-un text în vederea generării automate a rezumatului aceluși text (Ion et al., 2009b). Se observă însă că cele două probleme diferă în ceea ce privește numărul de clase care trebuie considerate și tipul entităților (pe de o parte paragrafe, iar de cealaltă, fraze) ce trebuie supuse procesului de clasificare. În cazul nostru, clasele pe care le-am considerat diferă ușor de cele furnizate de organizatorii competiției în specificațiile date. Astfel, clasele *reason* și *purpose* au fost unite, clasa obținută având denumirea de *reason-purpose*. Motivul constă în dificultatea dezambiguizării automate între cele două clase. Pentru a îmbunătăți precizia în faza de regăsire a paragrafelor relevante, am adăugat o altă clasă, etichetată *delete*, cu scopul de a elimina astfel, încă din faza de căutare, acele paragrafe care nu ar fi putut conține răspunsuri corecte pentru vreo întrebare. În această categorie intră paragrafele ce conțin titluri (e.g., „Articolul 1”), părți de tabele ale căror formatare nu s-a mai păstrat, semnături, etc.

Pentru faza de antrenare a clasificatorului am utilizat o colecție de aproximativ 800 de paragrafe etichetate manual cu următoarele etichete: *factoid*, *definition*, *procedure*, *reason-purpose* și *delete*, iar pentru a testa precizia, am folosit doar 89 de paragrafe. Metoda de clasificare la care am recurs este aceeași pe care am folosit-o și pentru clasificarea întrebărilor în anii precedenți: *principiul maximizării entropiei* (Ratnaparkhi, 1998). Trăsăturile pe care le-am luat în considerare s-au bazat pe cuvinte cheie, descriptori morfo-sintactici, punctuație, lungimea propoziției. Primele 5 cuvinte au fost considerate trăsături, ele având un puternic caracter de discriminare pentru clasele alese. Ca trăsături de tip morfo-sintactic, am luat în considerare descriptorul morfo-sintactic al verbului principal, o altă trăsătură fiind existența sau absența în cadrul propoziției a unui substantiv propriu. Alte trăsături sunt numărul de virgule din propoziție, numărul de ghilimele (indicând numărul de citate), semnul de punctuație cu care se termină propoziția. Din punct de vedere ortografic, o trăsătură importantă este numărul cuvintelor din propoziție care încep cu majusculă. Trăsăturile legate de lungime includ numărul propozițiilor în paragraf și lungimea paragrafului în cuvinte.

Precizia clasificatorului pe datele de test a fost de 94% (Ion et al., 2009b; Ștefănescu, 2010). Deși gradul de încredere statistică a evaluării preciziei este redus datorită numărului mic de 89 de paragrafe conținute de datele de test, rezultatele finale obținute folosind clasele asignate paragrafelor au fost mult îmbunătățite. În condițiile în care doar primele 50 de paragrafe întoarse de motorul de căutare au fost luate în considerare în etapele următoare ale fluxului de prelucrare, am constatat că, folosind clasificarea paragrafelor, în majoritatea cazurilor răspunsurile corecte s-au aflat în printre acestea.

Pentru clasificarea întrebărilor cele mai multe din sistemele actuale QA folosesc un modul specializat pentru a determina ce tip de răspuns ar trebui căutat în corpusurile avute la dispoziție. Desigur, clasificarea se poate realiza în mai multe feluri, cele mai simple metode folosind reguli sub forma unor expresii regulate. În cazul nostru, am apelat din nou la clasificatorul bazat pe maximizarea entropiei. Acesta a fost antrenat pentru a identifica 8 clase: *reason-purpose*, *procedure*, *definition*, *location*, *name*, *numeric*, *temporal* și *factoid*. Clasa *location* cuprinde întrebările care necesită ca răspuns o locație; clasa *name* conține întrebările care cer ca răspuns un nume de persoană, organizație sau numele unei entități (e.g., comisie, țară, etc.); clasa *temporal* cuprinde acele întrebări care cer drept răspuns o dată calendaristică sau un interval de timp; clasa *numeric* conține întrebările care cer drept răspuns un număr („Câți membri sunt în comisia de ...”), iar clasa *factoid* conține toate întrebările de tip factoid care nu sunt în clasele tocmai descrise. Trăsăturile pe care le-am folosit sunt următoarele: primul cuvânt de tip *WH* din întrebare, primul verb principal din întrebare, primul substantiv din întrebare, descriptorii morfo-sintactici ai tuturor substantivelor, verbelor, adjectivelor, adverbilor și numeralelor din întrebare, ordinea de apariție a primului verb și a primului substantiv în cadrul întrebării analizate. Desigur, extragerea acestor trăsături survine abia după preprocesarea întrebării.

Pentru antrenare, am plecat de la exemplele furnizate de organizatori și am construit 200 de întrebări (incluzând și acele exemple) pentru care am atașat manual clasa din care fac parte (Ștefănescu, 2010). Modelul a fost construit considerând doar trăsăturile care apăreau de cel puțin două ori în datele de antrenare, în final acesta având o acuratețe de 99% pe aceste date. Deși evaluarea performanțelor folosind datele de antrenare („*biased evaluation*”) trebuie evitată pentru aplicațiile de învățarea automată, ea ne-a permis totuși să ajungem la o selecție mai bună de trăsături. După cum era de așteptat, acuratețea a scăzut pentru cele 500 de întrebări furnizate odată cu startul competiției, însă la numai 97.2%. Scorul nu este atât de surprinzător pe cât pare, datorită diferențelor mari dintre trăsăturile caracteristice claselor considerate. De altfel, am ales clasele astfel încât clasificatorul să nu aibă probleme în identificarea lor corectă. Ne interesează în principal să nu avem erori de clasificare, chiar dacă rămânem cu clase nu foarte rafinate. Cu alte cuvinte, încercăm să restrângem cât mai mult spațiul de căutare fără însă a îngădui prea multe erori de clasificare (Ștefănescu, 2010).

5. Sistemul QA

Sistemul nostru este implementat ca un flux de prelucrare realizat peste arhitectura serviciilor și aplicațiilor web dezvoltate la ICIA. Sistemul este optimizat pentru a maximiza un scor de relevanță global $S(p)$ folosit la identificarea paragrafului cel mai plauzibil a constitui răspunsul adecvat la o întrebare adresată sistemului. Scorul global

$S(p)$ se calculează ca o combinație liniară a unor scoruri asigurate paragrafelor în funcție de criteriile de relevanță în raport cu o întrebare furnizată de un utilizator:

$$S(p) = \sum_i \lambda_i s_i, \quad \sum_i \lambda_i = 1 \quad (1)$$

unde s_i ($s_i \in [0,1]$) este unul din următoarele scoruri de relevanță:

- s_1 este 1 când clasificarea întrebării corespunde cu cea a paragrafului, în caz contrar, valoarea sa fiind 0. După cum am arătat, clasele modulului de clasificare a paragrafelor (5 la număr) diferă de cele ale modulului de clasificare a întrebărilor (8) dar între ele există o corespondență bine definită (Ion et al., 2009b);
- s_2 este un scor de similaritate lexicală bazat pe lanțuri lexicale, coeziune lexicală și identificarea perechii verb principal–argument; este scorul care încearcă să cuantifice gradul de similaritate lexicală dintre o întrebare și paragrafele întoarse de motorul de căutare ca fiind relevante pentru acea întrebare. Se calculează după următoarea metodă: pentru o întrebare Q și un paragraf candidat P, construim două liste conținând lemele ce corespund cuvintelor conținut din întrebare, respectiv din paragraf: LQ și LP. Numim lemele din lista LQ *cuvinte cheie*. Pe baza lor asignăm lui P un scor de relevanță. Astfel, s_2 se calculează ca un produs a trei alte scoruri ce caracterizează: (i) distanța semantică dintre lemele celor două liste (*DS*), (ii) coeziunea cuvintelor cheie în paragraf (*CC*) și (iii) identificarea în paragraf a unui posibil cuplu verb-argument extras din întrebare (*VA*). Aceste scoruri au fost descrise pe larg în (Ștefănescu, 2010) și nu vom mai insista asupra lor. Scorul final se calculează ca produsul celor trei scoruri:

$$s_2 = DS \times CC \times VA$$

- s_3 este un scor similar cu scorul BLEU (Papineni et al., 2002) care avantajează paragrafele în care cuvintele cheie ale întrebării apar în aceeași ordine ca în întrebare; ca și scorul precedent, se calculează considerând lemele cuvintelor conținut din întrebare, respectiv paragraf. Ideea implementării acestui scor (Ion et al., 2009b) are la bază observația că de foarte multe ori, în realitate, formularea răspunsului corespunzător unei întrebări conține o parte din acea întrebare. Principiul comparării n-gramelor a fost folosit în cazul de față pentru a evalua similaritatea dintre întrebare și paragrafele candidat.
- s_4 și s_5 sunt scorurile de relevanță pentru paragraf și document întoarse de motor.

4.1. Prelucrarea întrebării și alegerea răspunsului

După ce sistemul primește la intrare o întrebare, ea este trimisă serviciului web TTL pentru a fi preprocesată. Este apelat apoi serviciul web care se ocupă de clasificarea întrebărilor pentru a obține tipul de răspuns care trebuie căutat. În următorul pas, folosind informația adnotată după preprocesare, întrebarea este transformată în interogări într-un limbaj formal înțeles de motorul de căutare. Folosim 2 algoritmi pentru a genera două interogări diferite, ambele conținând ca termen de căutare clasa întrebării. Trebuie menționat că în faza de indexare au fost indexate odată cu paragrafele și clasele corespunzătoare acestora, pentru ca reducerea spațiului de căutare să se facă direct din faza de regăsire documentară.

Pentru fiecare din cele două interogări generate, motorul de căutare întoarce două liste L_1 și L_2 conținând fiecare 50 de paragrafe sortate după scorul de relevanță descris de

ecuația (1). Răspunsul întors de sistem este acel paragraf care se găsește atât în L_1 cât și în L_2 și care satisface:

$$\operatorname{argmin}_p(\operatorname{rang}_1(p) + \operatorname{rang}_2(p)), \operatorname{rang}_1(p) \leq K, \operatorname{rang}_2(p) \leq K, K \leq 50 \quad (2)$$

unde $\operatorname{rang}_1(p)$ este poziția paragrafului p în L_1 , iar $\operatorname{rang}_2(p)$, poziția paragrafului p în L_2 . Dacă nu există un paragraf care să satisfacă condițiile din (2), atunci sistemul semnalizează că nu poate găsi un răspuns la întrebarea dată, întorcând șirul de caractere NOA (*no answer* – nu există răspuns). În urma experimentelor, cele mai bune rezultate au fost obținute pentru $K=3$.

Sistemul dezvoltat este antrenabil, ponderile λ_i folosite în interpolarea liniară a scorurilor de relevanță fiind obținute cu ajutorul unei tehnici de optimizare similare cu metoda MERT. Metoda de antrenare pentru găsirea celui mai probabil răspuns a fost folosită pe cele 200 de întrebări utilizate și în cazul clasificatorului pentru întrebări și are următoarele etape: (i) rularea sistemului pentru cele 200 de întrebări și păstrarea primelor 50 de paragrafe întoarse de motorul de căutare pentru fiecare întrebare, ordinea paragrafelor fiind dată DOAR de scorurile de paragraf întoarse de motor (s_4); (ii) calcularea scorurilor s_i , $i = \overline{1,5}$, pentru fiecare din paragrafele extrase; (iii) pentru fiecare combinație de ponderi λ , cu $\sum_{i=1}^5 \lambda_i = 1$ și un pas de incrementare de 10^{-2} , se calculează scorul *Mean Reciprocal Rank* (MRR) (Radev et al., 2002) pentru întreg setul de 200 de întrebări, lista paragrafelor întoarse fiind sortată după ecuația (1); (iv) reținerea combinației de ponderi λ care maximizează scorul MRR.

Cele două moduri de a genera cereri în limbaj formal conduc la rezultate diferite ale motorului de căutare. Putem considera că avem de a face cu două sisteme ale căror rezultate sunt combinate pentru a obține răspunsul final. Astfel, cele două sisteme au fost individual optimizate, în ceea ce privește ponderile λ , fără a avea în vedere, în vreun moment, posibilitatea ca un sistem să întoarcă răspunsuri NOA. După antrenarea celor două sisteme, observând că sistemul de evaluare al competiției *ResPubliQA* favorizează răspunsurile NOA față de cele incorecte, am luat decizia de a introduce ecuația (2) pentru a păstra doar acele răspunsuri pentru care avem un grad de încredere ridicat. La momentul transmiterii rezultatelor, setul de ponderi λ utilizat nu ținea cont de tipul întrebărilor furnizate la intrare. Vom arăta în secțiunea de evaluare că putem obține rezultate îmbunătățite antrenând setul de ponderi pentru fiecare clasă de întrebări.

4.2. *Generarea de cereri formale pentru motorul de regăsire a documentelor*

Etapa de regăsire documentară este similară cu cea descrisă în (Ion et al., 2009a), folosind și de această dată motorul de căutare Lucene (Hatcher and Gospodnetić, 2004). Cuvintele documentelor au fost filtrate în funcție de descrierea lor morfo-sintactică astfel încât să rămână pentru indexare doar cuvintele conținut. Mai mult, cuvintele au fost normalizate la forma lor lemă. Eventualele erori de dezambiguizare morfo-lexicală sau lematizare au fost luate în calcul și pentru a compensa astfel de erori (mai ales în cazul unităților lexicale absente din dicționarul sistemului TTL) am indexat și forma ocurență pe aceeași poziție cu lema. În acest fel, un termen poate fi căutat atât ca lema, cât și în forma sa ocurență. Am procedat la construirea a doi indecși pentru colecția de documente: un index pentru paragrafe și unul pentru documente. Dată fiind o anumită cerere, motorul întoarce o listă alcătuită din paragrafele cele mai relevante în raport cu

acea cerere. Astfel, avem deja două scoruri calculate pentru un paragraf: scorul de relevanță pentru paragraf în indexul construit pe paragrafe (s_4) și cel de relevanță pentru documentul din care face parte paragraful în indexul construit pe documente (s_5).

Abilitatea sistemului nostru de a întoarce pentru unele întrebări mesajul că nici un răspuns nu a fost găsit se datorează combinării rezultatelor diferite obținute folosind cei doi algoritmi de generare a interogărilor formale către motorul de căutare.

ALGORITMUL DE GENERARE A CERERILOR FORMALE – VARIANTA 1 ia în considerare toate cuvintele conținut ale întrebării (substantive, verbe, adjective și adverbe) cu care construiește o disjuncție de termeni, ce sunt practic lemele acelor cuvinte. Singura condiție impusă asupra includerii unui termen în cerere este ca scorul TF-IDF (Salton and Buckley, 1988) al aceluși termen să fie peste un anumit prag (Ion et al., 2009b; Ștefănescu, 2010). În urma experimentelor, am stabilit ca valoarea acestui prag să fie 9. Folosirea acestui scor vizează eliminarea din interogare a termenilor comuni (cu frecvență de ocurență mare în documente) și deci, cu capacitate discriminatorie mică.

ALGORITMUL DE GENERARE A CERERILOR FORMALE – VARIANTA 2 folosește, asemeni variantei 1, informația obținută prin preprocesarea întrebării. Varianta precedentă a acestui algoritm, dezvoltată pentru căutările pe corpusul Wikipedia (Ion et al., 2009a) a fost extinsă și optimizată pentru noul corpus. Ca și în versiunea precedentă, algoritmul folosește verbele principale ale propoziției și grupurile nominale identificate în faza de preprocesare pentru a genera cererea formală către motorul de căutare. Noutatea constă în faptul că pentru fiecare grup nominal sunt construiți doi termeni ai cererii. (i) Primul este o expresie obținută prin concatenarea cu spațiu între ele a tuturor lemelor cuvintelor de tip conținut în ordinea apariției lor în grupul nominal. Setăm apoi 2 parametri pentru expresia formată. Primul, pe care îl vom numi *parametru de proximitate*, se referă în principal la numărul de cuvinte ce pot fi intercalate între termenii expresiei. Valoarea sa este egală cu 1 plus numărul cuvintelor funcționale ce se găsesc în grupul nominal. Al doilea parametru, pe care îl vom numi *parametru de amplificare*, definește importanța expresiei ca termen al interogării. Dacă acest parametru este setat la o valoare n , acest lucru înseamnă că identificarea ulterioară a termenului echivalează cu identificarea a n termeni obișnuiți. Am setat acest parametru la o valoare egală cu numărul unităților lexicale din expresia formată, pentru a favoriza paragrafele care conțin grupuri nominale identice cu cele din întrebare. (ii) Al doilea termen este o expresie Booleană formată doar din conjuncția lemelor corespunzătoare cuvintelor conținut din grupul nominal.

Observăm că primul termen este mai restrictiv pentru că impune o ordine de apariție și o anumită poziționare în documente a cuvintelor din întrebare. În schimb, pentru al doilea termen condițiile sunt relaxate. Urmărim ca atunci când cuvintele se găsesc în texte, dar nu strict în condițiile impuse de primul termen, să putem extrage totuși acele fragmente care le conțin. Ca și în cazul versiunii precedente, pentru fiecare verb principal din întrebare (cu excepția verbelor difuze semantic⁶) se generează un termen corespunzând lemei aceluși verb. Operatorul boolean folosit implicit de Lucene este SAU (OR) și astfel, cererile sunt disjuncții logice. Pe de altă parte, există și situația în care nici un paragraf nu conține o anumită expresie. În acest caz, cu ajutorul expresiilor conjunctive

⁶ Verbe cu putere discriminatorie redusă (a fi, a avea, a putea, etc.) întrucât utilizarea lor este comună în orice domeniu.

se încearcă identificarea celor care conțin cât mai multe cuvinte conținut ce apar și în întrebare.

Cele două cereri furnizate către motorul de căutare conduc la obținerea a două seturi de paragrafe ca rezultate. Putem vorbi astfel de două sisteme, chiar dacă, în foarte multe din componentele ce le compun, aceste sisteme sunt similare.

Una din cele mai utilizate practici pentru a îmbunătăți performanțele unui sistem QA este aceea de a îmbogăți termenii mono-cuvânt din interogările formale cu sinonime extrase din lexicoane sau ontologii lexicale precum Pinceton WordNet (Fellbaum, 1998). Conținând aproximativ 60.000 de mulțimi sinonimice, varianta românească a WordNet-ului (Tufiş et al., 2008a) este o excelentă resursă ce ar putea fi utilizată în acest scop. Cu toate acestea, datorită limbajului juridic specializat, caracteristic corpusului JRC-Acquis, am ales să nu facem uz de ontologia lexicală la nivelul generării interogărilor, ci într-o etapă ulterioară în care a fost necesară calcularea unui scor de similaritate lexicală dintre întrebare și paragrafele având probabilitate ridicată de a conține răspunsul corect.

5. Evaluare și concluzii

Evaluarea acurateței sistemelor înscrise în competiția ResPubliQA a fost făcută de organizatori folosind un program ce implementează formula (3):

$$A = \frac{1}{n} (n_{RC} + n_{NOA} \times \frac{n_{RC}}{n}) \quad (3)$$

unde n reprezintă numărul total de întrebări (500), n_{RC} este numărul de răspunsuri corecte, iar n_{NOA} reprezintă numărul întrebărilor la care am răspuns cu NOA. După cum se observă, formula de evaluare punctează suplimentar (termenul al doilea al formulei) capacitatea unui sistem de a discerne situațiile în care nu are certitudinea necesară de a da un răspuns corect și decide să emită un răspuns de tip „nu știu” (NOA=„no answer”). Ipoteza organizatorilor a fost că acuratețea în aceste situații ar fi aceeași cu acuratețea răspunsurilor corecte ($\frac{n_{RC}}{n}$). Cu alte cuvinte, atunci când sistemul decide că „nu știe” răspunsul corect, el este la fel de precis ca în cazurile în care furnizează un răspuns. Din analiza efectuată după primirea rezultatelor competiției reiese că în cazul sistemului nostru, precizia cu care răspunsurile greșite au fost evitate prin soluții NOA a fost chiar mai mare decât precizia răspunsurilor corecte.

Pentru evaluare, ICIA a trimis organizatorilor competiției două seturi de rezultate. Primul set, ICIA091RORO, conține rezultatele obținute prin combinarea ieșirilor celor două variante ale sistemului fără a folosi clasele întrebărilor ca termeni ai interogărilor formale, în timp ce al doilea set, ICIA092RORO, folosește aceste clase ca termeni. Utilizarea claselor întrebărilor ca termeni în interogările generate a condus la rezultate semnificativ îmbunătățite datorită faptului că în faza de indexare a paragrafelor am adăugat un câmp care să conțină clasa acestora.

Faza de antrenare desfășurată pe cele 200 de întrebări (v. formula (1)) a condus la obținerea următoarelor ponderi λ :

Tabelul 4: Ponderile pentru cele două variante de interogare ale sistemului

	λ_1	λ_2	λ_3	λ_4	λ_5
varianta 1 (interogări generate în funcție de scorul TF-IDF)	0,22	0,1	0,1	0,19	0,39
varianta 2 (interogări generate folosind structura de grupuri propoziționale a întrebărilor)	0,32	0,14	0,17	0,25	0,12

Fiecare din variantele sistemului a primit la intrare cele 500 de întrebări ale competiției, pentru ca apoi să întoarcă o lista a celor mai relevante 50 de paragrafe pentru fiecare întrebare, folosind parametrii din tabelul de mai sus. Dintre acestea, în cazul în care un paragraf a satisfăcut ecuația (2) acesta a fost furnizat ca răspuns. În caz contrar, răspunsul generat a fost NOA. Tabelul 2 conține evaluările oficiale ale rezultatelor trimise de ICIA: ICIA091RORO și ICIA092RORO. Scorul de acuratețe (formula 3) pentru sistemul ICIA092RORO a fost cel mai bun dintre toate cele 44 de seturi de răspunsuri (28 trimise de participanți, plus încă 16 scoruri de referință calculate de organizatori). Scorul de acuratețe pentru sistemul ICIA092RORO a fost al 4-lea.

Tabelul 5: Rezultatele oficiale ale ICIA la competiția ResPubliQA pe limba română

	ICIA091RORO)	ICIA092RORO
Întrebări la care a fost dat un răspuns	393	344
Întrebări la care NU fost dat un răspuns	107	156
Întrebări la care răspunsul a fost CORECT	237	260
Întrebări la care răspunsul a fost INCORECT	156	84
Întrebări la care răspunsul a fost NOA	107	156
scorul c@1	0,58	0,68

Interesant de menționat faptul că evaluarea cifrei de merit $M = \frac{A_{IR}}{A_{RD}}$ (v. (Peñas, et al., 2009) și Secțiunea 2) a pus în evidență mai obiectiv performanța sistemelor înscrise în concurs în raport cu limba de întrebare și răspuns. Cele două sisteme ale ICIA au obținut cele mai mari cifre de merit (ICIA092RORO=1,55, respectiv ICIA091RORO=1,32). Următoarele cifre de merit au fost obținute de sistemele de ÎR pentru limba italiană (1,24) și limba spaniolă (1,175). În raport cu scorul de acuratețe c@1, aceste două sisteme s-au clasat pe locurile 8 și respectiv 12. Aceste rezultate sunt foarte interesante și ele demonstrează un lucru mai puțin evident: pentru limba engleză câștigul de acuratețe adus de tehnologiile de limbaj natural față de tehnologiile standard de regăsire a informației este relativ mic (15%); acest lucru se explică prin faptul că limba engleză are o morfologie foarte simplă dar și prin faptul că cele mai avansate tehnici de regăsire documentară au fost dezvoltate, testate și optimizate pe documente scrise în limba engleză. Explicația se confirmă și prin faptul că analiza scorurilor de acuratețe ale sistemelor pentru diferite limbi cu morfologie semnificativă (germană, franceză, spaniolă, italiană, bulgară) au fost mai mici decât scorurile de acuratețe pentru limba engleză, deși unele dintre sistemele pentru limba engleză au fost elaborate de aceleași echipe care au implementat sistemele pentru limbile lor naționale. De asemenea, ierarhizarea în funcție de cifra de merit M, plasează cel mai bun sistem de ÎR pentru limba engleză⁷ pe locul al 6-lea, în timp ce sistemul pentru limba germană, cu al 17-lea scor c@1, are cifra de merit 1,158 care-l plasează pe poziția a 5-a.

⁷ Sistemul uned092 a obținut al doilea scor de acuratețe (c@1) 0,61, în timp ce cifra sa de merit (M) a fost de 1,151.

SISTEM ÎNTREBARE-RĂSPUNS ANTRENABIL PENTRU LIMBA ROMÂNĂ

Am testat și ipoteza conform căreia antrenarea ponderilor λ în funcție de clasa întrebării va conduce la îmbunătățirea rezultatelor. Am folosit încă odată cele 200 de întrebări pentru antrenare, partiționate în funcție de clasă, și astfel am obținut rezultatele din tabelul 3. Scorurile c@1 și M s-au îmbunătățit cu peste 3% față de cel mai bun rezultat precedent.

Tabelul 6: Ponderile rezultate după antrenarea pe fiecare clasă

		λ_1	λ_2	λ_3	λ_4	λ_5
varianta cu interogări generate în funcție de scorul TFIDF	Factoid	0,1	0	0,2	0,4	0,3
	Definition	0,2	0,15	0,05	0,15	0,45
	Reason-Purpose	0,1	0	0,15	0,3	0,45
	Procedure	0,1	0	0,15	0,15	0,6
varianta cu interogări generate folosind structura de grupuri propoziționale a întrebărilor	Factoid	0,15	0	0,3	0,3	0,25
	Definition	0,05	0,5	0,15	0,1	0,2
	Reason-Purpose	0,2	0	0,4	0,2	0,2
	Procedure	0,15	0,1	0,25	0,2	0,3

În final, trebuie să menționăm că sistemul nostru este disponibil on-line⁸, sub forma unei aplicații web ce apelează, pentru fiecare interogare, diferitele servicii web ce intră în arhitectura sistemului. În viitorul apropiat intenționăm să ne concentrăm eforturile în direcția dezvoltării unor capacități cross-linguale care să dea posibilitatea utilizatorilor să interogheze sistemul în limbaj natural fie în limba română, fie în engleză, iar sistemul să întoarcă rezultatele în oricare din aceste limbi. De altfel, secțiunea în limba engleză a corpusului JRC-Acquis a fost deja preprocesată și avem în vedere folosirea aplicațiilor de traducere automată dezvoltate la ICIA (Ceașu, 2009; Irimia, 2009) pentru a traduce fie întrebările în limbaj natural, fie cererile formale către motorul de căutare.

Mulțumiri. Activitatea de cercetare descrisă a fost sprijinită de CNMP – MECT prin proiectul național PNII SIR-RESDEC, D11007/18.09.2007.

Referințe bibliografice

- Ceașu, A. (2009). *Tehnici de traducere automată și aplicabilitatea lor limbii române ca limbă sursă*, București, România: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Fellbaum, C. (1998). *WordNet: An Electronic Lexical Database*, MIT Press.
- Hatcher, E. and Gospodnetić, O. (2004) *Lucene in Action*.
- Ion, R. (2007). *Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română*, București: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Ion, R., Ștefănescu, D., Ceașu, A. and Tufiș, D. (2009a). *RACAI's QA System at the Romanian-Romanian QA@CLEF2008 Main Task, Lecture Notes in Computer Science*, vol. 5706, September, p. 393–400.

⁸ <http://www2.racai.ro/sir-resdec/>

- Ion, R., Ștefănescu, D., Ceașu, A., Tufiș, D., Irimia, E. and Barbu-Mititelu, V. (2009b). *A Trainable Multi-factored QA System*, Proceedings of CLEF2009 Workshop, Corfu, Greece.
- Irimia, E. (2009). *Metode de traducere automată prin analogie. Aplicații pentru limbile română și engleză*, București: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Och, F.J. (2003). *Minimum Error Rate Training in statistical machine translation*, Proceedings of the 41st Annual Meeting on Association for Computational Linguistics, Sapporo, Japan, 160-167.
- Papineni, K., Roukos, S., Ward, T. and Zhu, W.-J. (2002). *BLEU: a method for automatic evaluation of machine translation*, Proceedings of the ACL-2002: 40th Annual meeting of the Association for Computational Linguistics, Philadelphia, USA, 311-318.
- Peñas, A., Forner, P., Sutcliffe, R., Rodrigo, Á., Forascu, C., Alegria, I., Giampiccolo, D., Moreau, N. and Osenova, P. (2009). *Overview of ResPubliQA 2009: Question Answering Evaluation over European Legislation*, Proceedings of the QA@CLEF workshop, Sept. 30 - Oct. 3.
- Pușcașu, G., Iftene, A., Pistol, I., Trandabăț, D., Tufiș, D., Ceașu, A., Ștefănescu, D., Ion, R., Orășan, C., Dornescu, I., Moruz, A. and Cristea, D. (2007). *Cross-Lingual Romanian to English Question Answering at CLEF 2006*, Lecture Notes in Computer Science, Available: ISBN: 978-3-540-74998-1.
- Radev, D.R., Qi, H., Wu, H. and Fan, W. (2002). *Evaluating Web-based Question Answering Systems*, Demo section, LREC 2002, Las Palmas, Spain.
- Ratnaparkhi, A. (1998). *Maximum Entropy Models for Natural Language Ambiguity Resolution*, Philadelphia, PA, USA: PhD Thesis, University of Pennsylvania.
- Salton, G. and Buckley, C. (1988). *Term-weighting approaches in automatic text retrieval*, Information Processing and Management, pp. 513-523.
- Steinberger, R., Poulliquen, B., Widiger, A., Ignat, C., Erjavec, T., Tufiș, D. and Varga, D. (2006). *The JRC-Acquis: A Multilingual Aligned Parallel Corpus with 20+ Languages*.
- Ștefănescu, D. (2010). *Extragere inteligentă de informații din corpusuri multilingve*, București: Teză de doctorat, Institutul de Cercetări pentru Inteligență Artificială, Academia Română.
- Ștefănescu, D. and Tufiș, D. (2006). *Aligning Multilingual Thesauri*, Proceedings of The 5th Language Resources and Evaluation Conference (LREC), Genoa, Italy.
- Tufiș, D., Ion, R., Bozianu, L., Ceașu, A. and Ștefănescu, D. (2008a). *Romanian Wordnet: Current State, New Applications and Prospects*, Proceedings of the 4th Global WordNet Conference (GWC-2008), Szeged, Hungary, 441-452.
- Tufiș, D., Ion, R., Ceașu, A. and Ștefănescu, D. (2008b). *RACAI's Linguistic Web Services*, Proceedings of the 6th Language Resources and Evaluation Conference (LREC 2008), Marrakech, Morocco.
- Tufiș, D., Ștefănescu, D., Ion, R. and Ceașu, A. (2008c). *RACAI's Question Answering System at QA@CLEF 2007*, Lecture Notes in Computer Science: Advances in Multilingual and Multimodal Information Retrieval (CLEF 2007), pp. 3284-3291.

UN SISTEM DE TRADUCERE AUTOMATĂ ROMÂNĂ - FRANCEZĂ

MIRABELA NAVLEA, AMALIA TODIRAȘCU

*LiLPa, Université de Strasbourg, Strasbourg – France, 22 rue René Descartes, BP
80010, 67084 Strasbourg cedex France;*

mirabela.navlea@yahoo.fr, todiras@unistra.fr

Rezumat

Această lucrare prezintă un proiect de cercetare având ca obiectiv adaptarea unui sistem de traducere automată pentru perechea de limbi română - franceză. Pentru traducere, se adoptă o metodă statistică, factorizată, care combină mai multe categorii de informații lingvistice (categorii lexicale, proprietăți morfologice și sintactice). Resursele pe care le folosește sistemul sunt corpusuri paralele adnotate și aliniate propozițional și lexical. Lucrarea prezintă corpusurile aliniate, precum și analiza lingvistică a erorilor de aliniere lexicală identificate în corpusurile paralele aliniate. Pe baza erorilor datorate diferențelor morfologice și sintactice dintre română și franceză, sunt propuse reguli euristice pentru a ameliora rezultatele alinierii lexicale.

1. Introducere

Obiectivul proiectului nostru este construirea de resurse lingvistice pentru un sistem de traducere automată statistică factorizată român - francez. Astfel, este studiată influența mai multor categorii de informații lingvistice (categorii lexicale, proprietăți morfo-sintactice) asociate unității lexicale, asupra calității traducerilor furnizate de un asemenea sistem.

Proiectul nostru adoptă metodologia propusă în *SEE-ERA.net* (Tufiș et al., 2008). Obiectivul acestui proiect era dezvoltarea de sisteme de traducere automată statistică factorizată pentru limbi slave și balcanice (bulgară, greacă, română, sârbă, slovenă), considerând engleza ca limbă sursă sau țintă. Aceste sisteme utilizează corpusuri paralele adnotate și aliniate propozițional și lexical.

Creșterea considerabilă a numărului documentelor disponibile în diferite limbi impune utilizarea de noi metode pentru manipularea bazelor de date voluminoase de documente multilingve, pentru crearea de situri web și aplicații multilingve, pentru ameliorarea motoarelor de căutare. Unele dintre aceste metode folosesc tehnici de traducere automată care trebuie să fie adaptate pentru diferite limbi, având în vedere faptul că majoritatea sistemelor de traducere automată consideră engleza ca limbă sursă sau țintă.

Engleza are o morfologie simplă, în raport cu alte limbi caracterizate de o morfologie flexionară complexă, cum este cazul limbilor slave și balcanice. Pentru acestea, construirea de resurse lingvistice (dicționare, gramatici) necesită timp și costuri materiale și umane importante. Cazul resurselor lingvistice românești este reprezentativ. Notăm existența mai multor dicționare, adnotatoare și corpusuri adnotate disponibile în limba română, dar majoritatea sistemelor de traducere automată propun româna și engleza (Marcu & Munteanu (2005); Irimia (2008); Ceaușu (2009); sistemul *Google*

*Translate*¹). La ora actuală, numai *Google Translate* tratează perechea de limbi română - franceză.

Mai mult, sistemele de traducere automată produc un număr important de erori de traducere datorate lipsei de resurse lingvistice performante pentru diferite perechi de limbi. Crearea de resurse lingvistice este o sarcină dificilă condiționată de numeroasele fenomene lingvistice care se produc într-o limbă dată (ambiguități, lipsa echivalențelor de traducere de la o limbă la alta, frazeologie). Grass (2009) identifică treisprezece tipuri de erori frecvente generate de sistemele de traducere automată: polisemia și omonimia, ambiguitatea (sintactică, referențială), termenii vagi (*fuzzy hedges*), expresiile și metaforele, neologismele, substantivele proprii, cuvintele de origine străină, împrumuturile și calcurile, separatorii, siglele și acronimele, transpoziția. Probleme specifice apar în cazul limbilor cu morfologie flexionară complexă, cum sunt româna și franceza. Astfel, numărul ridicat de forme flexionare ale cuvintelor mărește numărul de ipoteze de traducere. Pentru evitarea erorilor de acest fel, sistemele lingvistice de traducere automată se concentrează pe crearea de resurse lingvistice complexe, ca : gramatici, dicționare, baze terminologice sau de cunoștințe.

Dacă sistemele lingvistice (*Systran*²) obțin astfel rezultate de traducere performante, alte sisteme furnizează rezultate comparabile bazându-se pe tehnici statistice factorizate (*EuroMatrix*³, 2009), ce exploatează corpusuri paralele adnotate și aliniat. Sistemele statistice factorizate extind metodele bazate pe segmente de traducere (Koehn, Och & Marcu, 2003), ce utilizează doar forma de ocurență a cuvintelor, prin exploatarea mai multor factori lingvistici asociați unității lexicale, ca: leme, descrieri morfo-sintactice, informații sintactice etc. Așadar, sistemele statistice factorizate sunt modulare: diferiți factori lingvistici pot fi utilizați în procesul de traducere. Koehn și Hoang (2007) utilizează proprietățile morfo-sintactice, Avramidis și Koehn (2008) exploatează informația sintactică pentru a ameliora calitatea rezultatelor de traducere.

Astfel, pentru sistemul prezentat, se adoptă o metodă de traducere statistică factorizată, adaptând însă resursele utilizate pentru perechea de limbi studiate. Corpusurile paralele disponibile pentru română și franceză, aliniat propozițional, adnotate și lematizate vor fi aliniat la nivel lexical. Se va studia influența factorilor lingvistici utilizați asupra calității traducerilor română - franceză.

În secțiunea următoare, sunt prezentate arhitectura sistemului de traducere automată statistică factorizată și informațiile lingvistice care sunt exploatate pentru a construi modele de limbă și de traducere. Corpusurile utilizate sunt prezentate în secțiunea 3. Diferențele dintre română și franceză sunt discutate în secțiunea 4. Procesul de aliniere și erorile reperate în corpusul aliniat lexical sunt descrise în secțiunea 5. Concluziile și perspectivele acestui studiu sunt prezentate în secțiunea 6.

2. Metodologia adoptată

Proiectul nostru vizează dezvoltarea unui sistem de traducere automată statistică factorizată pentru română și franceză. Acest sistem a fost inițial implementat pentru

¹ <http://translate.google.com/>

² <http://www.systransoft.com/>

³ <http://www.euromatrix.net/>

română și engleză⁴ (Ceașu, 2009). Sistemul utilizează un corpus paralel adnotat aliniat propozițional și lexical, și aplică diferiți factori lingvistici: forme de ocurență ale cuvintelor, leme, etichete morfo-sintactice (setul de etichete MSD propus de proiectul Multext⁵ pentru franceză (Ide & Véronis, 1994) și română (Tufiș & Barbu, 1997)), grupuri nominale sau verbale nerecursive (chunk-uri), colocații.

Sistemul își dovedește eficiența (Ceașu, 2009), în principal, pentru traduceri din domeniul juridic - administrativ. Sistemul utilizează decodorul *MOSES* (Koehn et al., 2007), cu diferite configurații ale parametrilor lingvistici optimizați (leme, MSD), stabilite în funcție de sensul procesului de traducere.

Pentru a adapta *MOSES* la o nouă pereche de limbi, este necesară construirea unui model de limbă, utilizând un corpus monolingv adnotat în limba țintă, și a unui model de traducere factorizat, utilizând un corpus paralel adnotat și aliniat propozițional și lexical. Apoi, decodorul *MOSES* caută traducerea cea mai probabilă utilizând modelele de limbă și de traducere construite.

Cum sistemul nu a fost adaptat pentru română și franceză, studiul nostru vizează construirea resurselor lingvistice necesare pentru obținerea modelelor de limbă și de traducere pentru aceste limbi. Astfel, se iau în considerare următoarele etape:

- 1) constituirea corpusurilor paralele;
- 2) procesarea corpusurilor (segmentare lexicală, lematizare, adnotare morfo-sintactică și la nivelul grupurilor sintactice nerecursive);
- 3) alinierea propozițională și lexicală a corpusurilor paralele;
- 4) corectarea erorilor de aliniere, în urma unei analize detaliate a acestor erori. Aceste date sunt utilizate pentru a reantrena modulul de aliniere;
- 5) construirea modelelor de limbă în limba țintă, utilizând corpusuri monolingve adnotate, și a modelelor de traducere factorizate, utilizând corpusuri paralele adnotate și aliniate, pentru română și franceză;
- 6) configurarea sistemului cu factorii lingvistici cei mai relevanți (optimizare);
- 7) analiza lingvistică a erorilor de traducere, reconfigurarea sistemului prin reluarea etapei 5 și optimizarea sistemului.

În secțiunile următoare, sunt prezentate corpusurile monolingve și paralele utilizate. Sunt prezentate, de asemenea, etapele de procesare, procesul de aliniere lexicală și un ansamblu de erori de aliniere reperate în corpusul aliniat român - francez.

3. Corpusurile

Se utilizează un corpus paralel român - francez extras din *JRC-Acquis* (Steinberger et al., 2006). *JRC-Acquis* se bazează pe corpusul paralel multilingv *Acquis Communautaire* (legislația UE, 1950-prezent) și este disponibil pentru 231 de perechi de limbi obținute din 22 de limbi oficiale ale UE. *JRC-Acquis* este aliniat la nivel de paragraf și disponibil gratuit în format XML. În proiectul nostru, se folosește un subset

⁴ <http://www.racai.ro/webservices>

⁵ <http://aune.lpl.univ-aix.fr/projects/multext/>

de 228 174 de perechi de fraze aliniat 1:1 din *JRC-Acquis* (5 357 017 de cuvinte în română, 5 828 169 de cuvinte în franceză), selectate din ansamblul de documente comune în română și franceză.

Se folosește, de asemenea, un corpus paralel român - francez aliniat propozițional, extras din *DGT Translation Memory (DGT-TM)*⁶. Acest corpus se bazează tot pe *Acquis Communautaire*, dar are avantajul că majoritatea frazelor aliniat este corectată manual. *DGT-TM* în română și franceză conține 490 962 de perechi de fraze aliniat (9 142 291 de cuvinte în română și 9 953 360 de cuvinte în franceză). Acest corpus e disponibil gratuit în format TMX.

Cum *JRC-Acquis* și *DGT-TM* sunt corpusuri specializate juridic - administrativ, se utilizează și alte corpusuri paralele, constituite folosind resurse disponibile pe Internet, pentru a testa sistemul de traducere automată și pentru alte domenii (politică, aviație). Pentru colectarea corpusului paralel român - francez, au fost luate în considerare criterii ca: disponibilitatea textelor bilingve, fiabilitatea sursei textelor, calitatea traducerilor, precum și domeniul. Astfel, corpusul paralel a fost selectat și organizat manual ținându-se cont de criteriile menționate și specificându-se pentru fiecare text sursa indicată prin URL, autorul, data. Datele colectate au fost curățate prin eliminarea elementelor netextuale, ca: imagini, tabele, note de subsol etc. Pentru rezolvarea problemei lipsei diacriticelor pentru majoritatea textelor românești colectate de pe Internet, s-a folosit sistemul de recuperare a diacriticelor, *Diac+* (Tufiș & Ceaușu, 2008). Astfel, textele bilingve disponibile în română și franceză au fost colectate folosind mai multe surse web: documente de pe situl Parlamentului European (263 788 de cuvinte), siteurile companiilor aeriene românești (63 353 de cuvinte).

Pentru evaluarea sistemului nostru, se utilizează un corpus român - francez aliniat propozițional și lexical, constituit din 1000 de fraze (Todirașcu et al., 2008). Acesta a fost obținut printr-un proces de derivare (Tufiș & Koeva, 2007), folosindu-se două corpusuri paralele extrase din *JRC-Acquis* (englez - român și englez - francez) și aliniat lexical în mod automat. Corpusul aliniat lexical român - francez a fost corectat manual.

Pentru construirea modelelor de limbă, se utilizează următoarele corpusuri monolingve pentru limba română⁷:

- (i) partea românească a corpusului jurnalistic paralel *NAACL* (800 000 de cuvinte) (Martin, Mihalcea & Pedersen, 2005);
- (ii) corpusul *LT4eL* (alcătuit din manuale de utilizare, 600 000 de cuvinte) (Trandabăț et al., 2006);
- (iii) corpusul jurnalistic *RoCo* (7,5 milioane de cuvinte) (Tufiș & Irimia, 2006).

Pentru limba franceză, se folosesc corpusurile următoare:

- (i) un corpus juridic - administrativ selectat din *JRC-Acquis* (498 788 de cuvinte);
- (ii) un corpus jurnalistic selectat din *Le Monde* (1980-1988) (488 543 de cuvinte).

Corpusurile se procesează prin aplicarea tagger-ului *TTL*⁸ (Ion, 2007), disponibil pentru română și franceză ca serviciu web. *TTL* segmentează lexical, lematizează și adnotează

⁶ <http://langtech.jrc.it/DGT-TM.html>

⁷ Corpusuri disponibile la cerere de la autori.

textul cu descrieri morfo-sintactice și la nivelul grupurilor nominale, prepoziționale și verbale nerecursive (chunk-uri). Rezultatele furnizate sunt în format XCES (Figura 1) și conțin setul de etichete MSD propus de proiectul Multext pentru franceză (Ide & Véronis 1994) și română (Tufiş & Barbu, 1997).

```
<seg lang="fr"><s id="ttlfr.3">
<w lemma="voir" ana="Vmps-s">vu</w>
<w lemma="le" ana="Da-fs" chunk="Np#1">la</w>
<w lemma="proposition" ana="Ncfs" chunk="Np#1">proposition</w>
<w lemma="de" ana="Spd" chunk="Pp#1">de</w>
<w lemma="le" ana="Da-fs" chunk="Pp#1,Np#2">la</w>
<w lemma="commission" ana="Ncfs" chunk="Pp#1,Np#2">Commission
</w>
<c>.</c>
</s></seg>
```

Figura 1: Un exemplu de rezultate furnizate de TTL pentru franceză.

4. Construirea corpusurilor aliniate

După cum este menționat și în secțiunea anterioară, se utilizează corpusuri paralele adnotate și aliniate propozițional și lexical, pentru construirea modelelor de traducere factorizate. Calitatea corpusurilor aliniate este esențială pentru a obține rezultate de traducere performante. Cum nu dispunem de resurse pentru alinierea lexicală română - franceză, se utilizează corpusul aliniat propozițional (Ceașu et al., 2006), care este, de asemenea, lematizat, etichetat și adnotat la nivelul grupurilor nerecursive.

Se pregătește corpusul în formatul de intrare cerut de sistemul statistic *GIZA++*⁹ (Och & Ney, 2003), furnizându-se, de asemenea, lema și informația de categorie lexicală pentru dezambiguizarea morfologică a lemei. Se efectuează o aliniere lexicală bidirecțională, apoi se obține intersecția aliniamentelor.

Se filtrează lista echivalențelor de traducere astfel obținuți, cu o listă de cuvinte ce au ortografie asemănătoare și un sens comun (numite *cognates* în engleză) pentru română și franceză. Pentru identificarea acestor cuvinte, se folosește un algoritm care calculează cea mai lungă secvență de caractere comună pentru două cuvinte date. Dacă lungimea celei mai lungi secvențe comune de caractere este cel puțin 70% din lungimea celui mai scurt cuvânt, atunci perechea de cuvinte este considerată *cognat*.

Se utilizează lista filtrată a echivalențelor de traducere, ca punct de plecare pentru efectuarea alinierii lexicale.

Cum a fost propus în Tufiş et al. (2005), se aplică un ansamblu de euristici pentru realizarea alinierii lexicale: (i) definirea claselor de echivalență a categoriei lexicale (un substantiv poate fi tradus printr-un substantiv, verb sau adjectiv); (ii) alinierea substantivelor, adjectivelor, verbelor, adverbilor; (iii) alinierea grupurilor nominale, prepoziționale sau verbale nerecursive conținând echivalenți de traducere; (iv) alinierea elementelor conținute într-un astfel de grup, prin aplicarea de reguli euristice.

⁸ TTL = Tokenizing, Tagging and Lemmatizing free running texts.

⁹ <http://fjoch.com/GIZA++.html>

Pentru ameliorarea rezultatelor alinierii lexicale, au fost analizate erorile sistematice de aliniere. Astfel, a fost definit un ansamblu de reguli euristice contextuale pentru evitarea acestor erori. În secțiunea 5.2., sunt prezentate câteva dintre regulile definite.

5. Erorile de aliniere lexicală

S-a aliniat un corpus român - francez conținând 1000 de fraze extrase din *JRC-Acquis*, așa cum este descris în secțiunea anterioară. Au fost analizate apoi erorile sistematice de aliniere lexicală. Cele mai multe dintre acestea apar din cauza diferențelor morfo-sintactice dintre limbi, alte erori fiind specifice domeniului: cologații, termeni.

5.1. Comparație între română și franceză

În ciuda faptului că româna și franceza sunt amândouă limbi romanice, fiecare posedă caracteristici morfologice și sintactice specifice. Structurile sintactice sunt asemănătoare, dar caracteristicile morfo-sintactice prezintă diferențe importante. Astfel, substantivele și nominalele din română sunt caracterizate de caz (nominativ, acuzativ, genitiv, dativ, vocativ), marcat flexionar prin desinențe specifice și analitic, prin morfeme sau afixe proclitice (*GALR*¹⁰, 2005), și pot fi de genul neutru. Determinantul definit din română este întotdeauna enclitic și fuzionează cu substantivul și cu adjectivul antepus, în timp ce în franceză acesta este întotdeauna proclitic, formând cuvânt separat. În română, cliticile de acuzativ și de dativ se exprimă simultan cu complementul direct sau indirect. În franceză, utilizarea pronumelui clitic exclude realizarea de suprafață a complementului direct sau indirect ca grup nominal. Alte diferențe privesc structurile sintactice specifice fiecărei limbi (subordonatele relative), construcțiile verbale specifice (morfeme adiționale pentru anumite moduri sau timpuri în română, auxiliare specifice pentru verbe de mișcare în franceză), și morfemele diferite de la o limbă la alta (relația de posesie).

5.2. Analiza lingvistică a erorilor de aliniere lexicală

În această secțiune, sunt prezentate câteva dintre cele mai frecvente erori de aliniere lexicală reperate în corpusul studiat, extras din *JRC-Acquis*. Aceste erori sunt specifice domeniului juridic - administrativ al corpusului studiat, caracterizat de folosirea construcțiilor impersonale, a persoanei a-III-a singular sau plural etc. Pentru corectarea erorilor, a fost definit un ansamblu de 27 de reguli euristice morfo-sintactice contextuale. Aceste resurse se utilizează pentru efectuarea alinierii lexicale în interiorul grupurilor nerecursive (chunk-urilor). Se dau câteva exemple de erori și de reguli, acestea fiind explicate în (Navlea & Todirașcu, 2010).

Subordonatele relative. Una dintre erorile frecvente privește subordonatele relative. În subordonatele relative din română, complementul indirect este dublu exprimat, prin pronumele relativ *care* (dativ) și formele neaccentuate ale pronumelui personal *îi*, *-i*, *le*, *li*. În subordonatele relative din franceză, complementul indirect se exprimă cu ajutorul prepoziției *à* sau *auprès de* și al pronumelui relativ *lequel*. Din cauza acestei diferențe, în exemplul din Figura 2, pronumele personal *li* este nealiniat cu pronumele relativ *auxquelles* (formă contractată *à + lesquelles*).

¹⁰ Gramatica limbii române.

UN SISTEM DE TRADUCERE AUTOMATĂ ROMÂNĂ - FRANCEZĂ

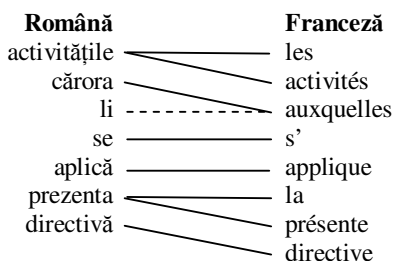


Figura 2: Eroare de aliniere lexicală privind subordonatele relative cu complement indirect.

Un exemplu de regulă euristică morfo-sintactică contextuală pentru evitarea acestui tip de erori, este dat în tabelul 1 de mai jos :

Tabel 1: Regulă euristică pentru subordonatele relative cu complement indirect

Română	Franceză
N + lema <i>care</i> (dativ) + <i>îil-illeli</i> + V	N + <i>àlauprès de</i> + lema <i>lequel</i> + V

Destinatarul. Erori de aliniere frecvente apar și în cazul exprimării destinatarului. În franceză, destinatarul este exprimat cu ajutorul prepoziției *à*, în timp ce în română acesta se exprimă prin dativ. Din cauza acestei diferențe, în exemplul din Figura 3, prepoziția *à* din franceză nu este aliniată cu substantivul în dativ din română.

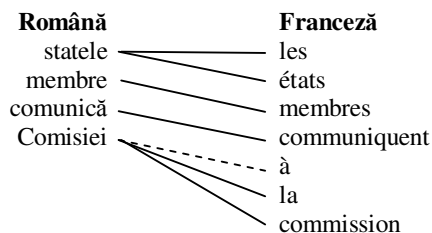


Figura 3: Eroare de aliniere lexicală privind exprimarea destinatarului.

Tabelul 2 de mai jos conține un exemplu de regulă euristică contextuală propusă pentru a evita această categorie de erori.

Tabel 2 : O regulă euristică pentru exprimarea destinatarului

Română	Franceză
V + N determinat definit (dativ)	V + <i>à</i> (<i>au, aux</i>) + determinant definit + N

Colocațiile. Se corectează, de asemenea, erorile de aliniere lexicală privind colocațiile. Colocațiile sunt expresii polilexicale, ale căror cuvinte sunt legate prin relații lexico-sintactice (Todirașcu et al., 2008). Aceste expresii nu sunt aliniat în bloc (Figura 4).

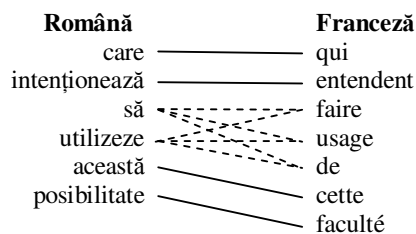


Figura 4: Eroare de aliniere lexicală privind cologațiile.

Se efectuează alinierea cologațiilor verbo-nominale prin utilizarea unui dicționar multilingv de cologații verbo-nominale, ce conține și proprietățile lor morfo-sintactice contextuale (Todirașcu et al., 2008). Cologațiile sunt utilizate astfel ca indici principali în procesul de aliniere lexicală. Cum dicționarul a fost completat cu date extrase din corpusul *JRC-Acquis*, și conține cele mai frecvente cologații verbo-nominale, existente în corpusurile francez, român, englez și german, acesta este eficient pentru alinierea cologațiilor verbo-nominale, însă nu rezolvă problema altor clase de cologații (nominale, adverbio-adjectivale etc.). Pentru acestea, se vor căuta alte strategii de aliniere.

6. Concluzii

Lucrarea de față prezintă un proiect în desfășurare ce urmărește construirea de resurse lingvistice pentru sisteme de traducere automată statistică factorizată, pentru două limbi cu morfologie flexionară complexă: româna și franceza. Sunt prezentate resursele utilizate pentru construirea corpusurilor aliniate lexical. Sunt analizate rezultatele modulului de aliniere lexicală și sunt propuse reguli euristice morfo-sintactice contextuale pentru ameliorarea acestor rezultate. Până în acest moment, au fost implementate regulile euristice și evaluarea rezultatelor nu este încă terminată. În viitor, se vor construi modele de traducere factorizate pentru română și franceză. Se vor compara rezultatele cu modele de traducere statistice pure și se vor evalua mai multe combinații de factori lingvistici, pentru găsirea parametrilor optimali în ceea ce privește perechea de limbi studiate.

Mulțumiri. Autorii mulțumesc doamnei Rada Mihalcea (University of Texas) pentru corpusul *NAACL*, d-lui Dan Cristea (Universitatea din Iași, România) pentru corpusul *LT4eL* și d-lui Dan Tufiș (Academia Română, București) pentru corpusul *RoCo*.

Referințe bibliografice

- Avramidis, E., Koehn, P. (2008). Enriching Morphologically Poor Languages for Statistical Machine Translation, In *Proceedings of ACL-08: HLT*, Columbus, June 2008, pp. 763-770.
- Ceașu, A., Ștefănescu, D., Tufiș, D. (2006). Acquis Communautaire Sentence Alignment using Support Vector Machines, In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)*, Genoa, May 2006, pp. 2134-2137.

- Ceașu, A. (2009). *Tehnici de traducere automată și aplicabilitatea lor limbii române ca limbă sursă*, Teză de doctorat, Academia Română, București, aprilie 2009, 123 p.
- Grass, T. (2009). A quoi sert encore la traduction automatique? *Les Cahiers du GEPE, Outils de traduction - outils du traducteur?*, n° 3, Strasbourg, 14 p.
- Guțu Romalo, V. (coord.) (2005). *Gramatica limbii române, Cuvântul*, vol. I, Ed. Academiei Române, București, 712 p.
- Ide, N., Véronis, J. (1994). Multext (multilingual tools and corpora). In *Proceedings of the 15th CoLing*, Kyoto, August 5-9, pp. 90-96.
- Ion, R. (2007). *Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română*, Teză de doctorat, Academia Română, București, mai 2007, .
- Irimia, E. (2008). Experimente de Traducere Automată Bazată pe Exemple, *Atelierul de Lucru Resurse Lingvistice Românești și Instrumente pentru Prelucrarea Limbii Române*, Iași, 19-21 novembre 2008, pp. 131-140.
- Koehn, Ph., Och, F. J., Marcu, D. (2003). Statistical Phrase-Based Translation. In *Proceedings of HLT-NAACL 2003*, Edmonton, May-June 2003, pp. 48-54.
- Koehn, Ph., Hoang H., (2007). Factored translation models, In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, Prague, June 2007, pp. 868-876.
- Koehn, Ph., Hoang, H., Birch, A., Callison-Burch, Ch., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, Ch., Zens, R., Dyer, Ch., Bojar, O., Constantin, A. Herbst, E. (2007). Moses: Open source toolkit for statistical machine translation. In *Proceedings of the ACL 2007*, Prague, June 2007, pp. 177-180.
- Marcu, D., Munteanu, D. S. (2005). Statistical Machine Translation: An English-Romanian Experiment, *EUROLAN 2005*.
- Martin J., Mihalcea R., Pedersen T. (2005). Word Alignment for Languages with Scarce Resources. In *Proceeding of the ACL2005 Workshop on Building and Using Parallel Corpora: Data-driven Machine Translation and Beyond*. June, 2005, Ann Arbor, Michigan, Association for Computational Linguistics, 65-74
- Navlea, M., Todirașcu, A. (2010). Linguistic Resources for Factored Phrase-Based Statistical Machine Translation Systems, In *Proceedings of LREC Workshop Exploitation of Multilingual Resources and Tools for Central and (South) Eastern European Languages*, Valletta, May 2010, pp. 41-48.
- Och, F. J., Ney, H. (2003). A Systematic Comparison of Various Statistical Alignment Models. *Computational Linguistics*, vol. 29, n° 1, March 2003, pp. 19-51.
- Steinberger, R., Pouliquen, B., Widiger, A., Ignat, C., Erjavec, T., Tufiș, D., Varga, D. (2006). The JRC-Acquis: A Multilingual Aligned Parallel Corpus with 20+ Languages, In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)*, Genoa, May 2006, pp. 2142-2147.
- Todirașcu, A., Heid U., Ștefănescu, D., Tufiș D., Gledhill C., Weller M., Rousselot F. (2008). Vers un dictionnaire de collocations multilingue. *Cahiers de Linguistique*, vol. 33, n° 1, Louvain, août 2008, p. 161-186.

- Trandabăț, D., Iftene, A., Pistol, I., Forăscu, C., Cristea, D. (2006). Resurse românești în cadrul proiectului LT4eL. În C. Forăscu, D. Tufiș, D. Cristea (eds.): Resurse lingvistice și instrumente pentru prelucrarea limbii române, Editura Universității „Al.I. Cuza” Iași, România, ISBN 978-973-703-208-9.
- Tufiș, D., Barbu, A., M. (1997). A Reversible and Reusable Morpho-Lexical Description of Romanian, In Dan Tufiș and Poul Andersen (eds.), *Recent Advances in Romanian Language Technology*, pp. 83-93, Editura Academiei Române, București, 1997. ISBN 973-27-0626-0.
- Tufiș, D., Ion, R., Ceaușu A., Ștefănescu, D. (2005). Combined Aligners. In *Proceedings of the ACL Workshop on Building and Using Parallel Texts: Data-Driven Machine Translation and Beyond*, pp. 107-110, Ann Arbor, USA, Association for Computational Linguistics. ISBN 978-973-703-208-9.
- Tufiș, D., Irimia, E. (2006). RoCo_News - A Hand Validated Journalistic Corpus of Romanian. In *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC 2006)*, pp. 869-872, Genoa, Italy, May 2006. ELRA - European Language Resources Association.
- Tufiș, D., Koeva, S. (2007). Ontology-supported Text Classification based on Cross-lingual Word Sense Disambiguation. In Masulli, F., Mitra, S., Pasi, G., eds., *Applications of Fuzzy Sets Theory. 7th International Workshop on Fuzzy Logic and Applications (WILF 2007)*, volume 4578 of Lecture Notes in Artificial Intelligence, September 2007, Springer - Verlag, pp. 447-455.
- Tufiș, D., Koeva, S., Erjavec, T., Gavrilidou, M., Krstev, C. (2008). Building Language Resources and Translation Models for Machine Translation focused on South Slavic and Balkan Languages. In Tadić (M.) et al., eds, In *Proceedings of the Sixth International Conference Formal Approaches to South Slavic and Balkan Languages (FASSBL 2008)*, Dubrovnik, September 25-28, pp. 145-152.
- Tufiș, D., Ceaușu, A. (2008). DIAC+: A Professional Diacritics Recovering System. In *Proceedings of the 6th Language Resources and Evaluation Conference - LREC 2008*, Marrakech, Morocco, May, 2008, ELRA, ISBN 2-9517408-4-0.

SERVICII WEB INTEROPERABILE ȘI MULTILINGUALE

RADU ION, ALEXANDRU CEAUȘU, DAN ȘTEFĂNESCU, DAN TUFIȘ

*Institutul de Cercetări pentru Inteligență Artificială, Academia Română
{radu, aceausu, danstef, tufis}@racai.ro*

Rezumat

Problemele interoperabilității instrumentelor și resurselor lingvistice sunt preocupări majore ale cercetării actuale în domeniul prelucrării limbajelor naturale. Cele mai noi tehnologii, bazate pe arhitecturi orientate pe servicii web, constituie un pas important în direcția asigurării interoperabilității. Atunci când serviciile web sunt independente de limbă sau facil adaptabile la diverse limbi criteriile de interoperabilitate și multilingualitate sunt în mare măsură satisfăcute. Orchestrarea diverselor servicii se face de obicei prin alegerea unui format comun al datelor de intrare/ieșire. În acest sens, au apărut deja o serie de platforme de integrare a resurselor și instrumentelor lingvistice dezvoltate și rezidente în locații diferite, implementate în limbaje de programare și sub sisteme de operare variate. În această lucrare prezentăm câteva dintre serviciile web ale Institutului de Cercetări pentru Inteligență Artificială (ICIA) adaptate la platforma de servicii web WebLicht. WebLicht reprezintă un mediu în care serviciile web înregistrate pot interacționa unele cu altele datorită conversiei parametrilor de intrare/ieșire la formatul TCF (engl. Text Corpus Format).

1. Introducere

Tehnologiile Web Service, un palier fundamental în filosofia noilor generații ale web-ului, înlesnesc utilizatorilor dezvoltarea de aplicații ce integrează diverse module, implementate de autori diferiți, în limbaje diferite, și chiar aflate pe alte mașini decât cea locală. Conceptul de arhitectură orientată spre servicii (SOA) a apărut dintr-o nevoie imperioasă de a structura și standardiza colecția eterogenă și vastă de instrumente și documente care există în spațiul virtual al intra- și inter-rețelelor informatice locale, regionale sau mondiale. Derivat din conceptul SOA, dar mai general, este conceptul de infrastructură de servicii web care oferă răspunsuri la construcția de fluxuri de prelucrare pe baza serviciilor web dar și la identificarea serviciilor relevante pentru o anumită aplicație, a datelor necesare diferitelor servicii, a publicării de resurse și instrumente de prelucrare a acestor resurse, a documentării lor etc. O astfel de infrastructură este în curs de construcție în cadrul proiectului european CLARIN¹.

CLARIN își propune să reunească aplicații și resurse din domeniul Prelucrării Automate a Limbajului Natural (PLN) și să le prezinte sub o formă în care nespecialiștii în prelucrarea limbajului natural (PLN) să le poată identifica și utiliza cât mai ușor în activitățile pe care le întreprind. Altfel spus, CLARIN urmărește realizarea unei infrastructuri europene a tehnologiilor limbajului natural (scris și vorbit) astfel încât acestea să fie *integrate* (aplicațiile și resursele se află în centre specializate care utilizează servere dedicate interconectate prin rețele de tip GRID), *interoperabile*

¹ <http://www.clarin.eu>

(resursele și serviciile vor fi descrise cu limbajele proprii Web-ului semantic (engl. Semantic Web) pentru a depăși barierele puse de diversitatea formatelor de intrare/ieșire), *stabile* (garantarea funcționării neîntrerupte și existența personalului de suport tehnic), *persistente* (garantarea funcționării continue pe o perioadă lungă de timp – cel puțin pe durata proiectului), *accesibile* (resursele și aplicațiile sunt accesibile online pe Web) și *extensibile* (întreg mediul trebuie să încorporeze ușor noi aplicații și resurse). CLARIN reunește și armonizează la nivel european o multitudine de proiecte naționale. Astfel de proiecte sunt, printre altele, proiectul german D-SPIN și proiectul românesc CLARIN-RO.

D-SPIN² (Hinrichs et al., 2008) a implementat o platformă de servicii web care să asigure interoperabilitatea unor aplicații de PLN dezvoltate independent. Metoda de interconectare a fost aceea a expunerii acestor aplicații ca servicii web de tip REST³ (și nu SOAP⁴ pentru că s-a considerat că volumul de lucru suplimentar necesar împachetării/despachetării formatului SOAP va introduce un timp suplimentar și inutil de adaptare pentru noile aplicații care doresc înregistrarea). Această platformă se numește **WebLicht**⁵ și a fost dezvoltată în colaborare de universitățile din Tübingen, Stuttgart, Leipzig și Berlin. În prezent conține peste 70 de aplicații de PLN (pentru limbile germană, engleză, franceză, română, spaniolă, italiană și finlandeză) expuse ca servicii web și care *sunt compatibile între ele ca formate de intrare/ieșire*. Acest lucru permite în mod evident posibilitatea extraordinară a compunerii de operații simple pentru obținerea unor rezultate care altfel ar fi fost imposibil de obținut – de exemplu, combinarea adnotării cu etichete morfosintactice a unui text din două surse (programe) diferite. Dintre operațiile ce se pot efectua asupra unui text putem menționa: segmentare lexicală, adnotare cu etichete morfo-sintactice, lematizare, analiză sintactică etc.

Formatul parametrilor de intrare/ieșire este unul XML care a fost definit special pentru a facilita integrarea aplicațiilor care prelucrează texte. Acesta se numește Text Corpus Format (TCF, (Schmid, 2009)) și reprezintă o stivă de adnotări care sunt produse de operații elementare de procesare a textelor (de exemplu, segmentarea la nivel de unitate lexicală este o adnotare necesară operației de etichetare morfo-sintactică). Astfel, fiecare serviciu web din WebLicht așteaptă la intrare un fișier XML care conține textul prelucrat până la un anumit nivel și întoarce același fișier XML la care adaugă un nou nivel care conține adnotările pe care le produce (niciunei operații nu-i este permis să șteargă vreun nivel de adnotare). Menționăm că platforma Weblicht este una dintre componentele ECD (European Clarin Demonstrator) prototipul funcțional al etapei de specificare (preparatory phase) a proiectului CLARIN, prototip a cărui implementare este coordonată de ICIA.

În cele ce urmează, vom descrie procesul de adaptare, în cadrul proiectului RO-CLARIN, a câtorva dintre serviciile web ale ICIA la platforma WebLicht. O serie de alte servicii web proprii sau în curs de implementare la ICIA vor fi adăugate în viitorul

² <http://weblicht.sfs.uni-tuebingen.de/englisch/index.shtml>

³ Representational State Transfer, un tip de comunicare de tip client-server.

⁴ Ajuns la versiunea 1.2, SOAP (<http://www.w3.org/TR/soap12-part0/>) este un protocol descris în XML cu care diverse aplicații distribuite pe calculatoare conectate în rețea, pot comunica prin schimb de mesaje împachetate în formatul SOAP.

⁵ <http://weblicht.sfs.uni-tuebingen.de:8080/WebLicht1/>

apropiat platformei WebLicht. De asemenea vom exemplifica avantaje ale adoptării de standarde pentru ușurința adaptării unor aplicații diverse (expuse ca servicii web) la rezolvarea unor probleme care altfel nu ar fi putut fi abordate. În acest sens, vom arăta cum anume putem combina rezultatele mai multor programe de etichetare morfo-sintactică operând pe aceeași segmentare lexicală cu metoda descrisă în (Tufiș, 1999) care opera cu același program de etichetare morfo-sintactică antrenat însă pe diverse corpusuri.

2. TTL în platforma WebLicht

TTL⁶ este un modul Perl care oferă următoarele adnotări ale textelor în limbile română, engleză și franceză: segmentare la nivel de frază și unitate lexicală, adnotare cu etichete morfo-sintactice (engl. POS tagging), lematizare, recunoașterea entităților denumite (engl. named entity recognition) și analiză sintactică de suprafață (engl. chunking). A ajuns la versiunea 7.9 și a fost implementat deja ca serviciu web (Tufiș et al., 2007).

Pentru adaptarea TTL la WebLicht a fost nevoie de următoarele modificări asupra implementării anterioare ca serviciu web:

- renunțarea la suportul pentru SOAP întrucât WebLicht specifică comunicarea de tip REST ca fiind standardul acceptat. În acest caz, un client care dorește să folosească serviciul web va trimite o cerere HTTP POST către serverul de web care găzduiește serviciul comunicându-i acestuia un șir de caractere codificat UTF-8 care reprezintă un document XML TCF bine format. Serviciul web procesează documentul și îl returnează clientului cu stratul de adnotări suplimentar pe care l-a produs.
- renunțarea la serverul de web implementat de pachetul Perl **SOAP::Lite** și adoptarea serverului de web **Apache**⁷ ca gazdă pentru TTL. În acest fel am rezolvat o problemă gravă a serverului de web din SOAP::Lite care nu accepta decât o singură conexiune la serviciul web la un moment dat și care de altfel, era și instabil. Cu Apache a trebuit să utilizăm modulul Fast CGI⁸ care permite serverului de web să încarce o *singură dată* o instanță a serviciului web (moment în care se încarcă toate resursele necesare – operație costisitoare ca timp de execuție) iar apoi să servească mai mulți clienți folosind această instanță. De asemenea, FCGI permite distribuția uniformă a încărcării pe instanțele serviciului web existente și un întreg management al pornirii/opririi instanțelor serviciului web.

Serviciul web TTL pentru platforma WebLicht expune câte un URL pentru fiecare operație elementară pentru fiecare limbă în parte (`{lang}` poate lua una din valorile engleză (`en`), română (`ro`) sau franceză (`fr`):

- segmentare la nivel de frază și apoi la nivel lexical: `http://ws2.racai.ro/TTL-{lang}-tokenizer;`

⁶ Abrevierea din engleză pentru „Tokenizing, Tagging and Lemmatizing free running texts”.

⁷ <http://httpd.apache.org/>

⁸ FCGI, <http://www.fastcgi.com/>

- adnotare cu etichete morfo-sintactice: <http://ws2.racai.ro/TTL-{lang}-postagger;>
- lematizare: <http://ws2.racai.ro/TTL-{lang}-lemmatizer;>
- analiza sintactică de suprafață: <http://ws2.racai.ro/TTL-{lang}-chunker;>

În afară de aceste operații, orice serviciu web din WebLicht necesită conversia textelor primare (care nu sunt adnotate în vreun fel) la formatul TCF de bază (fișierul XML care este preluat de prima operație care, de obicei, este segmentarea la nivel de frază/unitate lexicală). În prezent, WebLicht conține astfel de servicii web care convertesc texte UTF-8 și RTF la formatul TCF de bază exemplificat în figura 1.

```
<?xml version="1.0" encoding="UTF-8"?>
<D-Spin xmlns="http://www.dspin.de/data" version="0.3">
  <tns:MetaData xmlns:tns="http://www.dspin.de/data/metadata">
    <tns:source></tns:source>
  </tns:MetaData>
  <tns:TextCorpus xmlns:tns="http://www.dspin.de/data/textcorpus" lang="ro">
    <tns:text>Acesta este o propoziție de test.</tns:text>
  </tns:TextCorpus>
</D-Spin>
```

Figura 1: Formatul TCF pentru propoziția „Acesta este o propoziție de test.”

Figura 2 reprezintă o sesiune de lucru în WebLicht. Primul pas pe care trebuie să-l facă utilizatorul este să selecteze limba textului pentru care dorește procesarea. Apoi, alegând orice operație din coloana din partea stângă, sistemul va adăuga automat în lanțul de prelucrare toate operațiile care aduc textul din forma brută în formatul TCF necesar la intrarea operației selectate de utilizator. De exemplu, dacă am fi selectat operația de adnotare cu etichete morfo-sintactice, WebLicht ne-ar fi adăugat automat în lanțul de prelucrare operațiile de conversie la formatul TCF de bază (cel din figura 1) și segmentare la nivel de frază și unitate lexicală.

Lanțul de prelucrare final apare în chenarul „Selected Tools:” iar apăsarea butonului „Run” produce rezultatele vizibile în partea de jos a coloanei din partea dreaptă. Acestea se pot descărca de pe site în formatul XML TCF produs de lanțul de prelucrare.

Anterior aminteam de posibilitatea extraordinară de a combina diverse operații WebLicht pentru a obține rezultate care altfel (în absența acestei platforme) ar fi fost destul de dificil (sau imposibil în cazul în care aplicațiile nu erau disponibile) de obținut. Studiul de caz este combinarea rezultatelor a două adnotatoare cu etichete morfo-sintactice care rulează pe același text segmentat lexical și frazal în același fel pentru ambele programe. Tufiș (1999) discută posibilitatea obținerii unei adnotări morfo-sintactice mai bune în cazul în care același program de adnotare morfo-sintactică rulează cu modele de limbă diferite. WebLicht ne permite însă să obținem o adnotare mai bună (aplicând deci aceleași procedee ca în lucrarea citată – calculul matricelor de confuzie) folosind două adnotatoare diferite: în cazul nostru, TTL și TreeTagger pentru limba engleză. În figura 3 avem rezultatele celor două programe pentru propoziția „A very simple test sentence is the test bed for the combined classifiers model.”

SERVICII WEB LINGVISTICE ALE ICIA ÎN CADRUL PROIECTULUI CLARIN

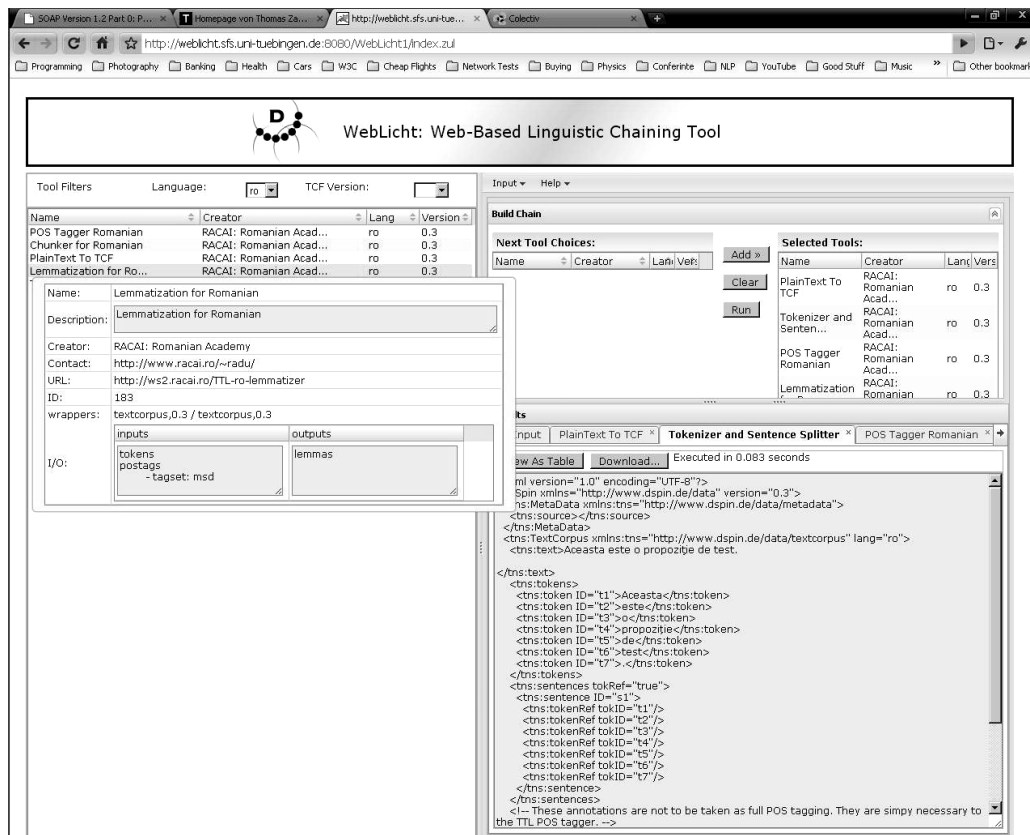


Figura 2: O sesiune de lucru în WebLicht

Pentru a reuși combinarea rezultatelor avem nevoie de corespondența etichetelor folosite (o altă aplicație extrem de utilă pe care WebLicht o face posibilă). TTL folosește mulțimea de etichete Multext-East (MSD) iar TreeTagger-ul folosește mulțimea de etichete Penn Treebank. Știm de exemplu că substantivele la singular se codifică cu „NN” în Penn Treebank și cu „Ncns” în MSD, substantivele la plural cu „NNS” respectiv „Ncnp”, adjectivele cu „JJ” respectiv „Afp”, ș.a.m.d. În exemplul considerat, ambele programe au găsit adnotarea corectă, dar, în cazul textelor mari, acest lucru este puțin probabil iar cu tehnica din (Tufiș, 1999) (după ce am fixat corespondența etichetelor), putem obține rezultate mai bune.

Tokenization	TreeTagger	TTL
A	DT	Ti-s
very	RB	Rsp
simple	JJ	Afp
test	NN	Ncns
sentence	NN	Ncns
is	VBZ	Vmip3s
the	DT	Dd

Tokenization	TreeTagger	TTL
test	NN	Ncns
bed	NN	Ncns
for	IN	Sp
the	DT	Dd
combined	JJ	Afp
classifiers	NNS	Ncnp
model	NN	Ncns
.	.	PERIOD

Figura 3: Rezultatele adnotatoarelor TTL și TreeTagger pe o propoziție de test

3. Servicii REST de procesare a textului bazate pe metode statistice de maximizare a entropiei

Pe lângă TTL, ICIA mai dispune de aplicații performante de prelucrare primară a textelor care au fost de asemenea adaptate ca servicii web SOAP așa cum se raportează în (Tufiș et al., 2007). Ca o noutate față de aceste servicii web de tip SOAP dezvoltate folosind platforma .Net Framework⁹ pentru engleză și română, noile servicii conforme cu specificațiile WebLicht (apel de tip REST și parametri de intrare/ieșire în format TCF) sunt disponibile și pentru limbile franceză și germană. Adaptarea la noile limbi a fost posibilă datorită flexibilității metodei de antrenare statistică, serviciul REST de adnotare morfo-sintactică¹⁰ folosind modulul de adnotare stratificată METT – Maximum Entropy Tiered Tagging (Ceașu, 2006), pe principiul maximizării entropiei, similar modelului ME al lui Ratnaparkhi (1988). Ca și TTL, tagger-ul METT utilizează pentru fiecare limbă un set de descrieri morfo-sintactice compatibil cu specificațiile Multex-East și un tagset redus pentru adnotarea stratificată (Tufiș & Dragomirescu, 2004).

Noile servicii web de procesare a textului care implementează METT pentru WebLicht sunt disponibile la adresa:

`http://www.racai.ro/RestWS/Service.svc/ws-{lang}-{webservice}`

unde {lang} este limba textului ce urmează a fi procesat, iar {webservice} este tipul de procesare ce urmează a fi invocat ca serviciu web. {lang} poate lua valorile română (ro), engleză (en), franceză (fr) și germană (de) iar {webservice} poate fi înlocuit cu unul din următoarele tipuri de prelucrări: segmentare lexicală (tokenizer), segmentare la nivel de frază (sentsplitter), adnotare morfo-sintactică (postagger) și lematizare (lemmatizer). Trebuie să notăm faptul că este obligatoriu ca în orice interogare HTTP la unul din aceste URL-uri, să specificăm în preambulul interogării

⁹ <http://www.microsoft.com/NET/>

¹⁰ <http://www.racai.ro/RestWS/Service.svc/ws-{lang}-postagger> unde „{lang}” poate lua valorile „en”, „ro”, „fr” și „de”.

tipul „application/xml” pentru fișierul TCF pe care îl vom trimite spre prelucrare (în antetul „Content-type”).

Pentru conversia în formatul TCF folosim serviciul web:

`http://www.racai.ro/RestWS/Service.svc/converter-{lang}`

unde {lang} este codul ISO de două caractere al limbii textului care urmează a fi procesat. Dacă limba nu este precizată, va fi invocat automat serviciul web pentru recunoașterea limbii (descriș în secțiunea 4).

În figura 4 este prezentat un exemplu de interogare a serviciilor web pentru conversia de la text la formatul TCF, segmentare lexicală și segmentare la nivel de frază, adnotare morfo-sintactică și lematizare. În continuare, folosind exemplul din figura 4, descriem parametrii de intrare și ieșire necesare fiecărui serviciu.

La interogarea serviciului web de conversie a unui text în formatul TCF (`http://www.racai.ro/RestWS/Service.svc/converter-ro`) rezultatul va fi un document XML similar celui din figura 1. Noul document XML va avea sub elementul „tns:TextCorpus” doar un element „tns:text” cu textul conținut în interogarea HTTP. În cazul acestui serviciu, parametrul de intrare este de tipul „text/plain” (specificat în „Content-type”) iar cel de ieșire este de tipul „application/xml”.

Pentru următoarea etapă de procesare – segmentarea lexicală – rezultatul serviciului de conversie este folosit ca argument pentru serviciul web `http://www.racai.ro/RestWS/Service.svc/ws-ro-tokenizer`. Rezultatul interogării HTTP este un document XML care are, pe lângă elementul „tns:text”, și un alt element „tns:tokens” care conține atomii lexicali din „tns:text”. Fiecare element „tns:token” conține indexul din textul inițial și lungimea atomului lexical.

În continuare, documentul XML rezultat al serviciului de segmentare lexicală este folosit ca parametru de intrare pentru serviciul de segmentare la nivel de frază (`http://www.racai.ro/RestWS/Service.svc/ws-ro-sentsplitter`). Acest serviciu adaugă elementul „tns:sentences” în care sunt grupate elemente „tns:sentence” reprezentând frazele formate din atomii lexicali din „tns:tokens”. Serviciul de segmentare frazală necesită prezența unui element „tns:tokens” în documentul XML inițial, deoarece fiecare frază conține referințe către atomii lexicali.

În final, rezultatul serviciului de segmentare la nivel de frază este folosit ca parametru de intrare pentru serviciul de adnotare morfo-sintactică stratificată (`http://www.racai.ro/RestWS/Service.svc/ws-ro-postagger`). Deși poate folosi ca parametru de intrare un document XML conținând doar segmentarea lexicală (cu element „tns:tokens”), este recomandat ca documentul XML ce urmează a fi prelucrat să conțină și segmentarea frazală (cu element „tns:sentences”), informația suplimentară contribuind la o viteză crescută a adnotării. Acest serviciu adaugă documentului XML inițial un nou element „tns:POSTags” în care sunt precizate descrierile morfo-sintactice ale atomilor lexicali.

Serviciul web de lematizare (`http://www.racai.ro/RestWS/Service.svc/ws-ro-lematizer`) necesită rezultatul tokenizării și al adnotării morfo-sintactice pentru a atribui o leamă

fiecărui atom lexical. Acest serviciu adaugă un element „tns:lemmas” documentului XML inițial.

```

<?xml version="1.0" encoding="utf-8" ?>
- <D-Spin version="3.0" xmlns="http://www.dspin.de/data"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema">
  <tns:MetaData xmlns:tns="http://www.dspin.de/data/metadata" />
- <tns:TextCorpus lang="ro"
  xmlns:tns="http://www.dspin.de/data/textcorpus">
  <tns:text>Acesta este un exemplu.</tns:text>
- <tns:tokens>
  <tns:token ID="t1" start="0" end="6">Acesta</tns:token>
  <tns:token ID="t2" start="7" end="11">este</tns:token>
  <tns:token ID="t3" start="12" end="14">un</tns:token>
  <tns:token ID="t4" start="15" end="22">exemplu</tns:token>
  <tns:token ID="t5" start="22" end="23">.</tns:token>
</tns:tokens>
- <tns:sentences>
  - <tns:sentence ID="s1">
    <tns:tokenRef tokID="t1" />
    <tns:tokenRef tokID="t2" />
    <tns:tokenRef tokID="t3" />
    <tns:tokenRef tokID="t4" />
    <tns:tokenRef tokID="t5" />
  </tns:sentence>
</tns:sentences>
- <tns:POSTags tagset="MSD">
  <tns:tag tokID="t1">Pd3msr</tns:tag>
  <tns:tag tokID="t2">Vmip3s</tns:tag>
  <tns:tag tokID="t3">Timsr</tns:tag>
  <tns:tag tokID="t4">Ncms-n</tns:tag>
  <tns:tag tokID="t5">PERIOD</tns:tag>
</tns:POSTags>
- <tns:lemmas>
  <tns:lemma tokID="t1">acesta</tns:lemma>
  <tns:lemma tokID="t2">fi</tns:lemma>
  <tns:lemma tokID="t3">un</tns:lemma>
  <tns:lemma tokID="t4">exemplu</tns:lemma>
  <tns:lemma tokID="t5">.</tns:lemma>
</tns:lemmas>
</tns:TextCorpus>
</D-Spin>

```

Figura 4: Exemplu de rezultat al interogării serviciilor web de identificare a limbii, de segmentare lexicală și frazală, de adnotare morfo-sintactică și de lematizare.

4. Identificarea limbii

Acest serviciu¹¹ asigură identificarea automată a limbii unui text scris într-una dintre cele 22 de limbi ale Uniunii Europene. Textul ar trebui să conțină un număr minim de 10-15 cuvinte (în principiu, o propoziție). Serviciul web de identificare a limbii

¹¹ <http://www.racai.ro/RestWS/Service.svc/converter>

funcționează ca un convertor din text în formatul TCF care poate recunoaște automat limba textului, adăugând un atribut „lang” elementului „tns:TextCorpus”.

Identificarea limbii este făcută folosind modele statistice pentru fiecare limbă și un modul de predicție. Predicția este realizată prin calcularea unui scor de similaritate al textului de intrare cu fiecare model de limbă în parte. Modelele de limbă se realizează pe baza ponderii pe care o au prefixele și sufixele cuvintelor în textele de antrenare disponibile pentru limba respectivă.

În experimentele realizate până acum, am utilizat texte de antrenament (cu mărimi ce au variat între 0,5 și 1,2 MB) pentru cele 22 limbi oficiale ale Acquis-ului Communautaire (Steinberger et al, 2006). Textele, fiind însă din domeniul juridic și având o structură mai aparte¹², nu sunt tocmai reprezentative pentru limbile luate în discuție. Cu toate acestea, am obținut rezultate excelente folosind pentru prefixe o lungime de trei caractere iar pentru sufixe de patru.

5. Concluzii

De la intrarea lor în funcțiune pe data de 19 februarie 2010, serviciile web ale ICIA din cadrul platformei WebLicht au avut 1389 de accesări pentru toate operațiile expuse pentru toate limbile. Dintre acestea cele mai multe accesări au fost de test (câte aprox. 40 de bytes de text pe cerere) dar au fost și cazuri în care s-a cerut procesarea a mai mult de 2KB de text. Suntem astfel siguri că efortul de integrare a aplicațiilor noastre în standardele promovate de CLARIN va fi de folos atât nouă cât și comunității CLARIN care crește într-un ritm alert.

Proiectul CLARIN deschide o nouă cale în abordarea cercetărilor în Prelucrarea Automată a Limbajului Natural și Lingvistică Computațională prin operațiunile de standardizare, colectare și diseminare a unor colecții impresionante de resurse și aplicații ale acestor domenii care altfel ar fi fost în marea majoritate a cazurilor, inaccesibile. Astfel, noi probleme de cercetare pot să apară sau unele mai vechi își pot găsi rezolvarea. În orice caz, adoptarea standardelor CLARIN va asigura accesul neîngrădit al celor neinițiați la tehnologiile limbajului pe de o parte și îmbunătățirea considerabilă a șanselor cercetătorilor de a găsi rapid soluții la problemele lor de cealaltă.

Mulțumiri. Activitatea de cercetare descrisă în această lucrare a fost sprijinită de proiectul european FP7 „CLARIN – Common Language Resources and Technology Infrastructure” (nr. de proiect 212230) finanțat de Comisia Europeană și de proiectul românesc PC7 „CLARIN–RO: Infrastructură pentru resurse lingvistice interoperabile pentru limba română” (nr. de proiect 16EU/06.04.2009) finanțat de ANCS.

¹² Textele juridice au adesea o structură formată din multe aliniate în care anumiți termeni se repetă de foarte multe ori afectând acoperirea lingvistică a modelelor de limbă.

Referințe bibliografice

- Ceașu, A. (2006). Maximum Entropy Tiered Tagging. Janneke Huitink & Sophia Katrenko (editors), *Proceedings of the Eleventh ESSLLI Student Session, ESSLLI 2006*, June 2006, Malaga, Spain, pp. 173—179.
- Hinrichs, E., Wittenburg, P., Lemnitzer, L., Geyken, A. (2008). D-SPIN - the German CLARIN initiative. In *CLARIN Newsletter #2*, 2008 (<http://www.clarin.eu/files/cnl02-web.pdf>).
- Kemps-Snider, M., Bel, N., Broeder, D. (2009). Proposal for a CLARIN Service CMDI Components. September, 2009, <http://www.clarin.eu/wp2/wg-26/wg-26-documents/cmdi-profile-for-web-services>
- Kemps-Snider, M., Bel, N. (2009). CLARIN Report on Web Services. March 2009 <http://www.clarin.eu/wp2/wg-26>
- Ratnaparkhi, A. (1998). Maximum Entropy Models for Natural Language Ambiguity Resolution. Ph.D. thesis, University of Pennsylvania, Philadelphia.
- Schmid, H. (2009). The technical details of the D-SPIN Architecture. Internal technical report, 2009 (<http://weblicht.sfs.uni-tuebingen.de/englisch/publikationen.shtml>).
- Steinberger, R., Pouliquen, B., Widiger, A., Ignat, C., Erjavec, T., Tufiș, D. (2006). The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. In *Proceedings of the 5th LREC Conference*, Genoa, Italy, 22-28 May, 2006, pp.2142-2147, ISBN 2-9517408-2-4, EAN 9782951740822.
- Tufiș, D. (1999). Tiered Tagging and Combined Language Models Classifiers. In Václav Matousek, Pavel Mautner, Jana Ocelíková, and Petr Sojka, editors, *Text, Speech and Dialogue (TSD 1999)*, Lecture Notes in Artificial Intelligence 1692, pp. 28—33. Springer Berlin / Heidelberg, January 1999. ISBN 978-3-540-66494-9.
- Tufiș, D., Dragomirescu, L. (2004). Tiered Tagging Revisited. In *Proceedings of the 4th LREC Conference*, Lisabona, 2004, pp. 39—42.
- Tufiș, D., Ion, R., Ceașu, A., Ștefănescu, D. (2007). Servicii Web lingvistice ale ICIA. În Ionuț Cristian Pistol, Dan Cristea, și Dan Tufiș, editori, *Lucrările atelierului RESURSE LINGVISTICE ȘI INSTRUMENTE PENTRU PRELUCRAREA LIMBII ROMÂNE*, pp. 61–68, Iași, România, 14–15 decembrie 2007. Editura Universității „Alexandru Ioan Cuza” Iași.

TIPARE HIPONIMICE PENTRU LIMBA ROMÂNĂ

VERGINICA BARBU MITITELU

Institutul de Cercetări pentru Inteligență Artificială, Academia Română

vergi@racai.ro

Rezumat

Hiponimia este o relație lexico-semantică foarte puțin studiată în lingvistica românească, însă ea oferă inginerilor o modalitate foarte eficientă de organizare a materialului lexical util în numeroasele aplicații dezvoltate, care presupun prelucrarea limbajului natural. Prezentăm în această lucrare două modalități de identificare a tiparelor hiponimice din limba română, rezultatele și evaluarea acestora; întrevădem, în final, utilități ale acestor tipare.

1. Introducere

Încă din antichitate au existat preocupări pentru ceea ce numim astăzi relații lexico-semantice. Contextul științific din a doua jumătate a secolului trecut a favorizat un reviriment al interesului pentru acest subiect, datorat conturării tot mai proeminente a unui nou domeniu de cercetare, lingvistica computațională, și a preocupărilor sale de prelucrare și generare a textelor în limbi naturale. Relațiile lexico-semantice intervin, la numeroase niveluri, în înțelegerea și producerea limbajului natural. În plus, ele reprezintă un factor semantic, care, alături de alții (frecvență, factorul stilistico-funcțional, etimologic și cel psihologic), contribuie la organizarea volumului uriaș de cuvinte ce alcătuiesc lexicul unei limbi (Bidu-Vrănceanu, Forăscu, 2005).

Hiponimia este relația paradigmatică de sens ce corespunde celei de incluziune între clase în logică (Lyons, 1977, Cruse, 1986). Cuvântul desemnând clasa mai cuprinzătoare din punct de vedere extensional se numește *hiperonim*, iar cel desemnând clasa cuprinsă este *hiponimul*. În lingvistica teoretică, această relație a fost definită, caracterizată (Cruse, 1986, Murphy, 2003), au fost propuse tipuri (Wierzbicka, 1984, Cruse, 1986, Kleiber, Tamba, 1990), teste de identificare (Cruse, 2004, Nyckees, 1998), a fost analizată prin raportare la relații asemănătoare precum meronimia și instanțierea (Kleiber, Tamba, 1990). În lingvistica computațională, hiponimia este privilegiată în ceea ce privește experimentele de extragere a relațiilor din texte, ei dedicându-i-se cele mai multe cercetări.

În această lucrare prezentăm două modalități de identificare a tiparelor hiponimice din limba română, utile pentru îmbogățirea wordnetului românesc cu noi sinseturi. După o descriere a stadiului actual al identificării hiponimiei în texte (secțiunea 2), descriem două modalități de identificare a tiparelor hiponimice din limba română (subsecțiunile 3.1 și 3.2), comparăm rezultatele celor două abordări (subsecțiunea 3.3), apoi prezentăm rezultatele evaluării acestor tipare pe un corpus de specialitate (secțiunea 4), urmate de prezentarea posibilelor aplicații ale acestui studiu (secțiunea 5) și de concluziile lucrării.

2. Stadiul cercetărilor

Studiile de lingvistică teoretică au adoptat preponderent o perspectivă paradigmatică asupra relațiilor lexico-semantice. Au fost identificate proprietățile semantice comune și enumerați termenii (atunci când aceștia erau în număr finit) alcătuind o paradigmă, cu ajutorul analizei semice sau componentiale. Relațiile au fost definite, caracterizate, clasificate după diverse criterii. Semanticienii au subliniat, însă, interdependența perspectivelor paradigmatică și sintagmatică în studiul semantic al cuvintelor (Bidu-Vrânceanu, Forăscu, 2005: 30).

Prin natura preocupărilor lor, lingvistica computațională și inteligența artificială sunt interesate, deopotrivă, și de planul sintagmatic. Având la dispoziție cantități uriașe de texte, cercetătorii pot extrage din ele diverse date, printre care și cuvinte aflate în diverse relații. Există două abordări principale ale extragerii relațiilor din texte: abordarea bazată pe tipare (engl. *pattern-based approach*) și cea bazată pe gruparea cuvintelor în funcție de sensul lor (engl. *clustering approach*). În cazul celei dintâi, există presupuziția că, la nivelul textului, pot fi identificate structuri lexico-sintactice cu grad înalt de specificitate în instanțierea unei perechi de cuvinte aflate într-o anumită relație semantică. Cele mai multe experimente s-au efectuat pentru hiponimie (Hearst, 1992, Alfonseca, Manandhar, 2001, Mann, 2002, Pașca, 2004, Oakes, 2005, Pantel, Pennacchiotti, 2006, Barbu Mititelu, 2008 etc.), meronimie (Berland, Charniak, 1999, Gîrju et al., 2006, Pantel, Pennacchiotti, 2006) și alte relații precum persoană – data nașterii, inventator – invenție, descoperitor – descoperire (Ravichandran, Hovy, 2002) etc. Abordarea bazată pe gruparea cuvintelor s-a concentrat asupra hiponimiei (Caraballo, 1999, Pantel, Ravichandran, 2004, Pantel, Pennacchiotti, 2006). Lucrarea de față se înscrie în primul tip de abordare.

3. Identificarea tiparelor hiponimice din limba română

În această lucrare folosim termenul *tipare hiponimice* pentru a ne referi la acele structuri lexico-sintactice care permit coocurența în text a unui hiperonim și a unuia dintre hiponimele sale, la scurtă distanță unul de altul, astfel: fie HIPERONIM *tipar hiponimic* HIPONIM, fie HIPONIM *tipar hiponimic* HIPERONIM. Am creat acest termen pentru a desemna ceea ce în literatura internațională de specialitate este numit *hyponymy/hyponymic pattern*.

Experimentele de identificare a tiparelor hiponimice din limba engleză sunt numeroase (vezi mai sus referințele). Din cunoștințele noastre, pentru limba română nu a mai efectuat nimeni asemenea experimente.

Prezentăm în continuare două modalități în care am identificat în texte românești tipare hiponimice.

3.1 Identificarea semi-automată a tiparelor hiponimice din limba română

Folosind un script Perl, am extras dintr-un corpus românesc (segmentat, lematizat, 881817 unități lexicale) propozițiile care conțin substantive aflate în relație de

hiponimie (directă sau indirectă¹), indiferent de distanța la care se află unul de celălalt în text. Pentru recunoașterea acestora, am utilizat wordnetul românesc, RoWN (Tufiş et al. 2008), care avea, la momentul respectiv, 46.269 de sinseturi. Nu am impus nicio restricție asupra distanței din arborele hiperonimic din wordnet. Propozițiile astfel extrase le-am grupat, automat, după asemănarea materialului lexical dintre hiponim și hiperonim: dacă în n propoziții apar un hiponim și un hiperonim (oricare ar fi ele) al său și între ele se află același grup de cuvinte, cu leme identice, în aceeași ordine, cele n propoziții, indiferent de hiponimul și hiperonimul conținute, sunt grupate împreună, ca reprezentând exemple ale aceluiași tipar hiponimic. Am evidențiat, în felul acesta, următoarele tipare (în care GN reprezintă un grup nominal, h=hiponim și H=hiperonim):

- GN(H) și **anume** GN(h): *Ea nu are nici_o bază de susținere, numai o idee, și anume o idee indestructibilă.*²
- GN(h) **fi un fel de** GN(H): *exprimare e un fel de pornografie*
- GN(H) **care avea fi** GN(h): *oameni care au fost participanți*
- GN(h) și (orice) **alt** GN(H): *bani și alte lucruri suspecte; taxele și orice alte venituri ale bugetului de stat*
- GN(h) și (tot) **celălalt** GN(H): *Legile și toate celelalte acte normative*
- GN(H), **mai ales** GN(h): *delicvenților de drept comun, mai ales gangsterilor*
- GN(h) **deveni** GN(H): *s-ar putea ca anumite prevederi să devină subiect de dispută*
- GN(H) **nu ca un** GN(h): *se poarte ca un om, nu ca un primar*
- GN(h) **fi considerat (ca)** GN(H): *televiziunile sunt considerate ca principale instrumente de luptă politică*
- GN(H) **sine numi** GN(h): *Acești oameni se numeau capitaliști*
- GN(H), **inclusiv** GN(h): *toți oamenii puterii, inclusiv miniștrii*
- GN(h) **fi un** GN(H): *Turcul e un om sărman, are întotdeauna mustață, face bani din piatră seacă și te întreabă dacă îți place Turcia*
- GN(h) **sau alt** GN(H): *bănci sau alte instituții de împrumut*
- GN(H), **adică** GN(h): *revin la „acest subiect sensibil”, adică la cazul Vântu - FNI*
- GN(H), **ci (și/doar) un** GN(h): *e un om cu idei, ci doar un animal*

Am testat aceste tipare pe un corpus jurnalistic de 900.000 de unități lexicale. Am extras automat din corpus acele propoziții în care apar structurile de mai sus. Manual, le-am selectat doar pe cele care îndeplineau și criteriul sintactic. Rezultatele sunt prezentate în Tabelul 1. Prin „număr de apariții” înțelegem numărul total de apariții ale structurii respective în corpus. Prin „număr de apariții relevante” desemnăm numărul de situații în care în pozițiile celor două grupuri nominale apar cuvinte aflate în relație de hiponimie.

¹ Hiponimia directă se stabilește între două noduri aflate în relație, din care unul este nodul-tată, iar celălalt nodul-fiu. Hiponimia indirectă se stabilește, datorită tranzitivității acestei relații, între noduri care nu sunt în relație unul cu celălalt, decât prin intermediul altui nod.

² Exemplul nu trebuie să ne surprindă, întrucât, în WordNet, cuvântul polisemantic *idee* apare în mai multe sinseturi, dintre care două sunt în relație de hiponimie: ex.: sinsetul {concept 1, idee1, noțiune2} cu glosa „idee generală care reflectă just realitatea” este hiponim al sinsetului {gând1, idee3} cu glosa „rezultatul procesului de gândire”.

Această relație poate fi deja înregistrată în RoWN, poate fi valabilă, deși, RoWN, incomplet fiind, nu o conține, sau poate fi o hiponimie strict contextuală, care nu necesită înregistrare în rețeaua semantică. Am efectuat manual această etapă de evaluare constrânși de două circumstanțe: incompletitudinea RoWN, pe de o parte, și dorința de a nu pierde din vedere cazurile de hiponimie stabilită strict contextual, care nu ar fi fost recunoscute automat. În plus, întrucât corpusul cu care am lucrat nu este parsat, evaluarea manuală a permis și luarea în considerare a acelor cazuri în care determinanți ai hiponimului sau hiperonimului se intercalează în structură, neafectând caracterul acestuia: de exemplu, „revin la „acest **subiect** sensibil”, adică la **cazul** Vântu - FNI”, unde hiperonimul *subiect* este urmat de un determinant (adjectivul *sensibil*), iar hiponimul *cazul* este precedat de prepoziția *la*.

Precizia tiparelor am calculat-o ca raport procentual între numărul relevant de ocurențe ale unui tipar și numărul său total de ocurențe. În cazul acestei analize, precizia este egală cu acoperirea (engl. *recall*), deci reprezintă chiar acuratețea acestor tipare.

Tabel 1: Rezultatele testării tiparelor hiponimice românești

Nr. crt.	Tipar	Număr de apariții	Număr de apariții relevante	Acuratețe (%)
1.	GN și orice alt GN	2	2	100
2.	GN și celălalt GN	4	4	100
3.	GN mai_ales GN	1	1	100
4.	GN fi considerat GN	3	3	100
5.	GN sine numi GN	6	6	100
6.	GN fi un GN	36	36	100
7.	GN sau alt GN	2	2	100
8.	GN ci (și/doar) un GN	3	3	100
9.	GN deveni GN	15	14	93,3
10.	GN și anume GN	11	10	90,1
11.	GN și alt GN	7	6	85,7
12.	GN inclusiv GN	31	23	74,2
13.	GN adică GN	8	5	62,5
14.	GN nu ca un GN	1	0	0

Menționăm că două tipare nu s-au regăsit în corpusul folosit pentru evaluare: *GN fi un fel de GN*, *GN care avea fi GN*.

3.2 Traducerea tiparelor englezești și regăsirea lor în corpusuri românești

În (Barbu Mititelu, 2008) am prezentat un experiment în care am identificat tipare hiponimice englezești. Redăm, în Tabelul 2, tiparele hiponimice englezești cele mai relevante și acuratețea lor (NP înseamnă *noun phrase* „grup nominal”).

Tabel 2: Tipare hiponimice englezești și acuratețea lor

Nr. crt.	Tipar	Acuratețe (%)
1.	NP other than NP	100
2.	NP especially NP	100

TIPARE HIPONIMICE PENTRU LIMBA ROMÂNĂ

Nr. crt.	Tipar	Acuratețe (%)
3.	NP principally NP	100
4.	NP usually NP	100
5.	NP such as NP	99,2
6.	NP in particular NP	92,3
7.	NP e(.)g(.)NP	91,4
8.	NP become NP	91
9.	NP another NP	87
10.	NP notably NP	86,8
11.	NP particularly NP	84,6
12.	NP except NP	84,6
13.	NP called NP	81,5
14.	NP like NP	81,3
15.	NP including NP	80,6
16.	NP mainly NP	75
17.	NP mostly NP	70,8
18.	NP i.e. NP	65

Menționăm că în afară de tiparele incluse în tabel, am mai identificat și altele, pentru care, însă, nu am putut stabili acuratețea, întrucât ele nu s-au regăsit în corpusul de testare. Este vorba despre tiparele: NP *be another* NP, NP *namely* NP, NP *and other* NP, NP *or other* NP, NP *a form of* NP, NP *or another* NP, NP *and similar* NP, NP *or similar* NP, NP *not least* NP, NP *but not* NP, NP *a kind of* NP, NP *like other* NP, NP *in common with other* NP, NP *and sometimes other* NP, NP *and many other* NP, NP *and in other* NP, NP *or any other* NP, NP *which be* NP, NP *for example* NP, NP *that is* NP, NP *apart from* NP, NP *even* NP, NP *be* NP, NP *for instance* NP, NP *as* NP, NP *either* NP, NP *as well as* NP. Am tradus aceste tipare și am verificat apariția lor într-un corpus românesc de 900000 de unități lexicale. Metoda de evaluare este identică cu cea descrisă mai sus, pentru tiparele românești identificate semiautomat. Rezultatele sunt înscrise în Tabelul 3.

Tabel 3: Tipare hiponimice românești obținute prin traducerea celor englezești și ocurența lor în corpus

Nr. crt.	Tipar	Număr de apariții	Număr de apariții relevante	Acuratețe (%)
1.	GN de exemplu GN	1	1	100
2.	GN ca de pildă GN	1	1	100
3.	GN de pildă GN	2	2	100
4.	GN cum ar fi GN	7	7	100
5.	GN, mai puțin GN	1	1	100
6.	GN de obicei GN	1	1	100
7.	GN altul/alta/alții/alte decât GN	2	2	100
8.	GN sau alt/altă/alți/alte GN	9	9	100
9.	GN sau orice alt/altă/alți/alte GN	2	2	100
10.	GN și anume GN	2	2	100
11.	GN numit GN	51	50	98

Nr. crt.	Tipar	Număr de apariții	Număr de apariții relevante	Acuratețe (%)
12.	GN deveni GN	114	101	88,6
13.	GN mai ales GN	8	7	87,5
14.	GN și alt/altă/alți/alte GN	53	44	83
15.	GN adică GN	30	19	63,3
16.	GN cu excepția GN	16	10	62,5
17.	GN care fi GN	8	5	62,5
18.	GN în special GN	5	3	60
19.	GN inclusiv GN	53	29	54,7
20.	GN afară de GN	11	5	45,5
21.	GN precum GN	6	2	33,3
22.	GN în afară de GN	9	2	22,2
23.	GN chiar și GN	29	6	20,7
24.	GN un fel de GN	27	5	18,5
25.	GN alt/altă/alți/alte GN	8	1	12,5
26.	GN ca GN	700	40	5,7
27.	GN până și GN	11	0	0
28.	GN dar nu GN	2	0	0

3.3 Comparație între tiparele hiponimice românești identificate prin cele două metode diferite

Diferența cantitativă între tiparele traduse și cele identificate își găsește trei explicații:

- diferența cantitativă între cele două wordneturi utilizate: wordnetul englezesc PWN 2.1 conține 117.597 de sinseturi, în vreme ce versiunea de RoWN folosită de noi conține 46.269 de sinseturi;
- diferența cantitativă între corpusurile utilizate: fragmentele din British National Corpus pe care le-am folosit pentru identificarea tiparelor englezești însumează 1863MB, iar cele românești 25,6 MB;
- atunci când am evaluat ocurențele tiparelor traduse, am acceptat ca relevante și acele situații în care se creau hiponimii contextuale, adică hiponimii neînregistrate în wordnet, dar stabilite între cuvinte care, în respectivul context, funcționează ca hiponime (vezi exemplul 1), precum și pe acelea în care hiponimia rezultă în urma decodării metonimice a contextului³ (vezi exemplul 2):

(1) ... să comercializeze chiloți, târlici, oale și alte comedii...

(2) Pentru cei care trec prin guvern, parlament și alte funcții publice, fenomenul de combatere a sărăciei se încheie cu o inexplicabilă îmbogățire.

³ Hearst (1992) notează și ea existența unor hiponimii rezultate din metonimie, dependente de context sau de perspectiva din care sunt privite lucrurile. Având în vedere scopul cu care extrage ea hiponime din text, adică acela de a îmbogăți automat WordNetul englezesc, problema nu este deloc trivială și necesită o cântărire a exemplelor înainte de a le introduce în WordNet.

Comparând mulțimea tiparelor românești identificate prin traducerea celor englezești (3.2) cu mulțimea celor identificate semiautomat (3.1), constatăm următoarele:

- intersecția celor două mulțimi, reprezentată de contextele identificate prin ambele metode de lucru, cuprinde: *sau alt*, *și anume*, *sine numi*, *deveni*, *mai_ales*, *și alt*, *adică*, *inclusiv*. Aceasta înseamnă că aproape 29% din tiparele traduse au fost identificate și semiautomat și că 57% din cele regăsite semiautomat au fost identificate și prin traducerea tiparelor englezești;
- aceste tipare comune celor două mulțimi sunt, așa cum era de așteptat, dintre cele cu acuratețe mare; în general, peste 80%, excepție făcând *adică* și *inclusiv*. Calculând o medie a preciziilor finale pentru fiecare dintre acestea, se pare că tiparul *sau alt* are precizie maximă; celelalte au o precizie medie de peste 90%: *numit/sine numi* 99%, *și anume* 95%, *mai_ales* 93,8%, *deveni* 91%; urmează apoi *și alt* cu 84,4%, *inclusiv* cu 64% și *adică* cu aproape 63%. În termeni lingvistici, cu cât precizia este mai mică, cu atât tiparul respectiv manifestă o „polisemie sintactică” mai mare, adică apare în mai multe structuri sintactice, cu valori diferite.
- tiparele comune nu sunt neapărat și cele mai frecvente în limbă, lucru iarăși așteptat. Precizia mare a acestora implică un polisemantism sintactic redus, deci o frecvență scăzută.

Dacă facem o comparație între tiparele englezești cu precizie maximă și cele românești cu precizie maximă, constatăm că echivalenții a trei dintre cele englezești se regăsesc printre tiparele cu precizie maximă din limba română: *other than* – *altul decât*, *especially* – *mai_ales*, *usually* – *de obicei*.

Din numărul total de tipare englezești, 70% au ca echivalent în limba română tipare hiponimice. Invers, din numărul total de tipare românești, 66,(6)% au ca echivalent în limba engleză tipare hiponimice. După cum se observă, procentele sunt foarte apropiate, ceea ce dovedește faptul că limbi diferite tind să permită coocurența în contexte similare a cuvintelor aflate în relație de hiponimie.

4. Evaluarea tiparelor hiponimice românești pe un corpus specializat

Am testat tiparele hiponimice din limba română pe un subcorpus al corpusului OPUS (<http://www.let.rug.nl/~tiedemann/OPUS/>), adică pe documentele EMEA (European Medicines Agency), cuprinzând 11.914.802 de unități lexicale și semne de punctuație. Acest corpus are două caracteristici care influențează rezultatele experimentului: abundența expresiilor repetate și vocabularul specializat (Tiedemann, 2009). Metoda evaluării este aceeași cu cea descrisă mai sus pentru evaluarea tiparelor identificate prin cele două metode diferite (vezi 3.1. și 3.2). Rezultatele sunt cuprinse în Tabelul 4.

Tabel 4. Testarea tiparelor hiponimice pe un corpus specializat

Tipar	Acuratețe (%)
GN <i>chiar și</i> GN	100
GN <i>de obicei</i> GN	100

Tipar	Acuratețe (%)
GN, <i>ci (și/doar)</i> GN	100
GN <i>în special</i> GN	96.88
GN <i>precum</i> GN	94.83
GN <i>cum ar fi</i> GN	93.75
GN (<i>în</i>) <i>afară de</i> GN	92.11
GN <i>și (orice) alt</i> GN	90.1
GN <i>fi un</i> GN	87.98
GN <i>sau alt</i> GN	86.96
GN <i>mai ales</i> GN	85.71
GN <i>alt decât</i> GN	85.71
GN <i>sine numi</i> GN	84
GN <i>inclusiv</i> GN	83.51
GN <i>de exemplu</i> GN	79.57
GN <i>fi considerat</i> GN	79.17
GN <i>care fi</i> GN	74.12
GN, <i>adică</i> GN	66.66
GN <i>cu excepția</i> GN	54.55
GN <i>și (tot) celălalt</i> GN	54.29

Comparând aceste rezultate cu cele obținute pe baza corpusurilor jurnalistice, remarcăm că, în majoritatea cazurilor, acuratețea scade. Sunt și câteva tipare care dovedesc o acuratețe mai mare în cazul textelor specifice unui domeniu: GN, *adică* GN, GN *care fi* GN, GN *în special* GN, GN *inclusiv* GN, GN *precum* GN, GN (*în*) *afară de* GN. Două dintre tipare au acuratețe maximă în ambele tipuri de texte: jurnalistice și specializate: GN *de obicei* GN și GN, *ci (și/doar)* GN.

5. Aplicații ale tiparelor hiponimice

Wordnetul românesc este o resursă lingvistică extrem de utilă aplicațiilor din lingvistica computațională. Calitatea și cantitatea sinseturilor conținute contribuie la calitatea aplicațiilor care folosesc wordnetul. În aceste condiții, îmbogățirea lui este una dintre prioritățile noastre. Dezvoltat manual până acum, i se pot adăuga sinseturi identificate automat cu ajutorul tiparelor hiponimice prezentate mai sus, după ce un specialist lingvist le adaugă o definiție și le verifică completitudinea.

Tiparele hiponimice pot fi utilizate pentru îmbogățirea RoWN, pe de o parte, atât cu hiponime, cât și cu instanțe, iar pe de altă parte, atât cu cuvinte din vocabularul general, cât și cu termeni din diverse domenii.

În cercetarea prezentată aici, nu am distins între hiponimie și instanțiere, întrucât ele erau tratate identic în versiunea de wordnet folosită. Acestea sunt, însă, relații semantice distincte, iar ultima versiune de wordnet englezesc, Princeton WordNet 3.0, le marchează diferit. Rămâne de cercetat în ce măsură tiparele identificate astfel își modifică precizia dacă facem această distincție și cât de precise sunt în identificarea instanțelor. Includerea instanțelor într-o resursă lingvistică precum wordnetul își dovedește cu prisosință utilitatea în sarcini de identificare a entităților numite (engl. *named entities recognition*).

În ceea ce privește cuvintele din limba comună, marea majoritate a acestora au fost deja incluse în RoWN. Interesul îl poate reprezenta îmbogățirea sa cu termeni din diverse domenii. Experimentul de testare a tiparelor hiponimice pe un corpus specializat aduce date relevante în această privință.

Din perspectiva semanticii lexicale, experimente ca cel prezentat aici pot aduce completări, exemplificări și rafinări analizei preponderent paradigmatică a hiponimiei.

6. Concluzii

Lucrarea de față se înscrie pe linia cercetărilor de extragere a relațiilor semantice din corpusuri cu ajutorul tiparelor lexico-sintactice. Am prezentat două modalități de identificare a tiparelor hiponimice din limba română. Una dintre metode este semiautomată, cealaltă presupune traducerea unor tipare hiponimice englezești și verificarea echivalentelor acestora în texte românești. Am calculat acuratețea tiparelor atât într-un corpus jurnalistic, cât și într-unul dintr-un domeniu.

Experimentul nu este deloc lipsit de utilitate: tiparele hiponimice pot fi folosite pentru identificarea în corpusuri a hiponimelor de introdus în wordnetul românesc, și chiar și a instanțelor.

Mulțumiri. Autoarea este recunoscătoare Ministerului Educației, Cercetării, Tineretului și Sportului, finanțator al proiectului SIR-RESDEC, în cadrul căruia s-a derulat cercetarea prezentată în acest articol.

Referințe bibliografice

- Alfonseca, E., Manandhar, S. (2001). Improving an Ontology Refinement Method with Hyponymy Patterns. *Third International Conference on Language Resources and Evaluation*, Las Palmas, 235-239.
- Barbu Mititelu, V. (2008). Hyponymy Patterns. Semi-automatic Extraction, Evaluation and Inter-lingual Comparison. *Text, Speech and Dialogue*, Springer (P. Sojka, A. Horak, I. Kopecek, P. Karel, eds.), 37-44.
- Berland, M., Charniak, E. (1999). Finding parts in very large corpora. *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, 57-64.
- Bidu-Vrăncianu, A., Forăscu, N. (2005) *Limba română contemporană. Lexicul*, București, Humanitas Educațional.
- Caraballo, S. (1999) Automatic construction of a hypernym-labeled noun hierarchy from text. *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, 120-126.
- Cruse, D. A. (1986). *Lexical Semantics*, Cambridge, CUP.
- Cruse, D. A. (2004). *Meaning in Language. An Introduction to Semantics and Pragmatics*, ediția a doua, Oxford, OUP.
- Gîrju, R., Badulescu, A., Moldovan, D. (2006). Automatic Discovery of Part-Whole Relations. *Computational Linguistics*, 32 (1): 82-135.

- Hearst, M. A. (1992). Automated Acquisition of Hyponyms from Large text Corpora. *Proceedings in the Fourteenth International Conference on Computational Linguistics*, Nantes.
- Kleiber, G., Tamba, I. (1990). L'hyponymie revisit e: inclusion et hi erarchie. *Langages* 98: L'hyponymie et l'hyperonymie, Larousse.
- Lyons, J. (1977). *Semantics*, vol. 1, Cambridge University Press.
- Mann, G. S. (2002). Fine-Grained Proper Noun Ontologies for Question Answering. *COLING-02 on SEMANET: building and using semantic networks*, 1-7.
- Murphy, M. L. (2003). *Semantic Relations and the Lexicon*, Cambridge, CUP.
- Nyckees, V. (1998). *La s emantique*, Paris, Belin.
- Oakes, M. P. (2005). Using Hearst's Rules for the Automatic Acquisition of Hyponyms for Mining a Pharmaceutical Corpus. *Proceedings of the Workshop Text Mining Research, Practice and Opportunities*, Borovets, 63-67.
- Pantel, P., Ravichandran, D. (2004). Automatically labeling semantic classes. *Proceedings of HLT/NAACL-04*, Boston, 321-328.
- Pantel, P., Pennacchiotti, M. (2006). Espresso: Leveraging Generic Patterns for Automatically Harvesting Semantic Relations. *Proceedings of Conference on Computational Linguistics / Association for Computational Linguistics (COLING/ACL-06)*, Sydney, 113-120.
- Pa ca, M. (2004). Acquisition of Categorized Named Entities for Web Search. *Proceedings of CIKM'04*, Washington.
- Ravichandran, D., Hovy, E. (2002). Learning surface text patterns for a question answering system. *Proceedings of ACL-2002*, Philadelphia, 41-47.
- Tiedemann, J. (2009). News from OPUS – A collection of multilingual parallel corpora with tools and interfaces. *Recent Advances in Natural Language Processing: Selected Papers from RANLP 2007*, John Benjamins, 237-248.
- Tufi , D., Ion, R., Bozianu, L., Ceau u, A.,  tef nescu, D. (2008). Romanian Wordnet: Current State, New Applications and Prospects. *Proceedings of 4th Global WordNet Conference, GWC-2008*, Szeged (Attila Tanacs, Dora Csendes, Veronika Vincze, Christiane Fellbaum, Piek Vossen, eds.), 441-452.
- Wierzbicka, A. (1984). Apples are not a „kind of fruit”. *American Ethnologist* 11, 313-28.

MECANISMELE GENERATIVE ALE MORFOLOGIEI DERIVAȚIONALE

PETIC MIRCEA

Institutul de Matematică și Informatică, Academia de Științe a Moldovei
mirsha@math.md

Rezumat

În articol sunt studiate problemele elaborării unui generator de derivate. Punctul de pornire este un lexicon de cuvinte derivate, care conține nu doar reprezentarea grafică a derivatelor, ci și a morfemelor lor constituente. Aceasta a permis studierea și formularea unor reguli care ar genera cuvinte derivate conform unor restricții la nivel de reprezentare grafică.

Cuvinte-cheie: generare automată a derivatelor, dicționar de derivate, prefix, sufix.

1. Introducere

Dicționarele moderne se confruntă cu unele deficiențe, devenind astfel obiect de cercetare pentru lexicografi. Având în vedere că dicționarele sunt permanent completate cu intrări noi, grație dezvoltării limbii, sarcina elaborării unui vocabular complet rămâne una practic imposibilă. Prin urmare, completarea automată și/sau semiautomată a resurselor lingvistice cu cuvinte generate automat în baza celor deja existente prin mijloace exclusiv interne, în particular cu ajutorul derivatelor cu prefixe și sufixe, reprezintă o sursă semnificativă de îmbogățire a vocabularului.

Pentru automatizarea procesului de derivare este necesar:

- ◆ să stabilim regulile, care pot fi aplicate cuvintelor de bază pentru obținerea unor cuvinte derivate noi;
- ◆ să stabilim condițiile, în care pot fi aplicate aceste reguli;
- ◆ să elaborăm și să aplicăm un mecanism de validare a acestor reguli, având în vedere că restricțiile menționate ulterior nu garantează complet corectitudinea cuvintelor obținute.

Pentru soluționarea problemelor sus menționate este necesar de efectuat un studiu preliminar al procesului de derivare. În calitate de suport lexicografic în acest scop au fost utilizate câteva surse și anume varianta electronică a dicționarului de derivate (Constantinescu, 2008), www.dexonline.ro și RRTLN¹ (*Resurse Reutilizabile ale Tehnologiei Limbajului Natural*).

Scopul acestui articol este de a studia particularităților cuvintelor derivate și de a stabili unele mecanisme generative ale morfologiei derivaționale.

În acest scop, inițial, este descrisă varianta electronică a dicționarului de derivate, evidențiindu-se astfel o caracterizare a dicționarului prin prisma unor date statistice în

¹ Lexiconul se conține pe site-ul <http://imi201.math.md/elrr/>

ce privește derivatele și constituentele lor. În continuare, este prezentat cazul particular de derivare semianalizabilă, care suscită mai multe semne de întrebare în ce privește modul de formare al lor. Ulterior, s-a studiat situația derivării regulate, care presupune, la rândul său, generarea derivatelor prin schimbarea genului cuvântului, generarea diminutivelor și a augmentativelor. Un compartiment aparte îl constituie procesul de proiectare a derivării cu prefixe cu o ulterioară derivare cu sufixe. În final sunt prezentate unele particularități ale generării automate a derivatelor cu prefixele *im-/in-*.

2. *Lexiconul cu derivate*

Lexiconul reprezintă varianta electronică a dicționarului de derivate, elaborat de S. Constantinescu (Constantinescu, 2008) și conține doar reprezentarea grafică a derivatelor și morfemelor constituente, fără nici o informație despre partea de vorbire a derivatelor și temelor lor. Pentru o procesare mai simplă a intrărilor lexiconului, a fost elaborată o expresie regulată care reprezintă următoarea structură a derivatelor:

$$\text{derivat} = (+\text{morfem}) * .\text{morfem}(-\text{morfem}) *$$

unde *+morfem* reprezintă un prefix, *.morfem* este tema și *-morfem* este un sufix. Un exemplu de intrări în dicționar este:

antistatal=+anti.stat-al
reprogramabil=+re.programa-bil

Au fost elaborați algoritmi și apoi dezvoltate programele (Petic, 2010) care stabilesc caracteristicile statistice ale lexiconului (Tabel 1).

Tabel 1: Caracteristicile statistice ale lexiconului

Caracteristica	Numărul
Derivate	15300
Rădăcini	6800
Prefixe	42
Sufixe	433

În momentul procesării lexiconului, prin extragerea grupurilor de derivate cu aceleași afixe, s-a constatat că există un grup restrâns de prefixe și sufixe care formează marea majoritate a derivatelor (Petic, 2010).

Astfel, 12 prefixe din cele 42 formează 88,2% din toate derivatele cu prefixe, înregistrate în acest lexicon. Derivatele care sunt formate din următoarele prefixe sunt cele mai numeroase: *ne-*, *re-*, *în-*, *des-*, *pre-*, *anti-*, *auto-*, *sub-*, *dez-*, *supra-*, *de-*, și *îm*².

Pe de altă parte, din 433 de sufixe înregistrate în lexicon 52 formează 87,7% din numărul tuturor derivatelor cu sufixe. Cele mai numeroase derivate s-au dovedit a fi, în ordine descrescătoare, următoarele sufixe: *-re*, *-tor*, *-toare*, *-eală*, *-ie*, *-ătoare*, *-iza*, *-oasă*, *-ar*, *-ător*, *-ească*, *-os*, *-aș*, *-esc*, *-tură*, *-iță*, *-ist*, *-uță*, *-el*, *-i*, *-ui*, *-ătură*, *-ește*, *-ism*, *-a*, *-ărie*, *-ică*, *-ime*, *-itate*, *-ioară*, *-ișor*, *-ișoară*, *-ic*, *-uleț*, *-că*, *-ean*, *-iș*, *-easă*, *-bil*, *-uț*, *-at*, *-oaică*, *-ușor*, *-an*, *-oi*, *-uliț*, *-iu*, *-enie*, *-istă*, *-al*, și *-ea*.

² Derivatele sunt enumerate în ordine descendentă.

E de menționat faptul că nu de la toate temele se formează un număr proporțional de derivate (Petic, 2010). Astfel, sunt cuvinte la care se înregistrează un număr maxim de derivate și anume: **bun** (32 de derivate), **alb** (25 de derivate), **șarpe** (22 de derivate), **roată** (22 de derivate), **om** (20 de derivate). În lexicon, un număr foarte mare de cuvinte (3657) apar cu un singur derivat. O ilustrație a acestei situații este prezentată în Fig.1.

3. Particularitățile derivatelor semianalizabile

După Iorgu Iordan (Iordan, 1970), derivatele în limba română pot fi grupate în: analizabile, semianalizabile și neanalizabile.

În derivatele *analizabile* se recunoaște atât prefixul, cât și cuvântul de bază. În derivatele *semianalizabile* se recunoaște numai prefixul, prin opoziție cu alte derivate sau cuvinte compuse cu o temă comună, inexistentă ca și cuvânt independent (de exemplu: *deschis* – *închis*). Unele derivate au devenit cu timpul *neanalizabile*, fie din cauza evoluției fonetice, morfologice sau semantice, fie din cauza dispariției cuvântului de bază la care se făcea raportarea.

Depistarea derivatelor semianalizabile pare a fi o problemă mai complicată decât cea a celor analizabile, deoarece temele care au stat la baza derivărilor lor nu mai sunt cunoscute. Astfel, pentru a studia acest proces, este necesar de a avea la dispoziție cuvintele derivate semianalizabile marcate corespunzător. Prin urmare, este util lexiconul, care conține cuvintele derivate împreună cu structura lor morfemică (Constantinescu, 2008). Analizând conținutul lexiconului menționat, s-a constatat că el include astfel de derivate, stabilindu-se, că, în mare parte, derivatele semianalizabile se referă la prefixele *des-/în-*. S-a realizat extragerea automată a acestor cuvinte, identificându-se 57 de derivate de acest tip. Așadar, unele merită atenție pentru a pătrunde în esența mecanismului de derivare.

O caracteristică pentru derivatele semianalizabile este utilizarea derivării regresive, adică înlăturarea unui afix și nu adăugarea acestuia, un exemplu elocvent fiind cuvântul *crucișătură* care s-a format prin derivarea regresivă de la verbul *încruși* și afixarea cu sufixul *-ătură*. Corectitudinea procesului de formare a acestui cuvânt este confirmată prin intrarea în dicționarul electronic de derivate:

crucișătură=. (în) cruciși-ătură

Deci, pentru acest caz se observă că intrarea din lexicon descrie modul de formare a derivatului. Verbul *încruși* a derivat regresiv, prin înlăturarea prefixului *în-*, ca apoi să fie adăugat sufixul *-ătură*. Intrarea din lexicon a marcat derivarea regresivă prin prefixul luat între paranteze rotunde. În resursa www.dexonline.ro, aceeași informație este prezentată după cum urmează:

[în]cruși + suf. -ătură

ceea ce reprezintă o marcare similară, cu cea de mai sus, a modului de formare a cuvântului derivat *crucișătură*. Dincolo de acest cuvânt, mai sunt și alte derivate de acest tip, de exemplu, *cingătoare*, *chiondoreală*, *fundătură*, *greunatic*, *lăuntric*, *notătoare*, etc.

Un alt exemplu de derivate semianalizabile sunt derivatele cu prefixele *des-* (*dez-*). Din exemplele ulterioare, se desprind următoarele cuvinte prefixate semianalizabil:

despăduri=+des. (îm)păduri

dezvălui=+dez. (în)vălui

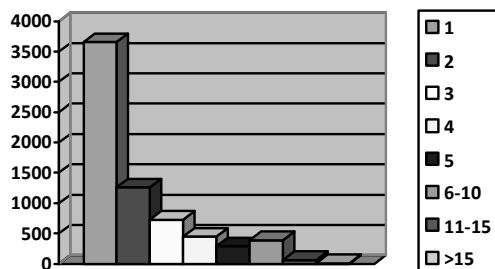


Figura 1: Datele statistice cu referire la numărul derivatelor cu temele din lexicon.

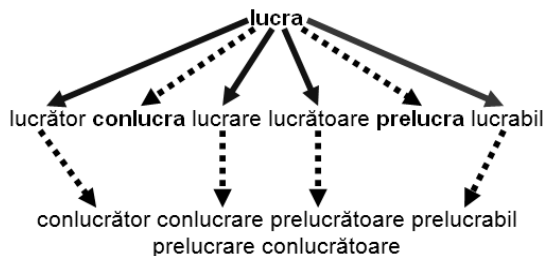


Figura 2: Proiectarea derivatelor în cazul cuvântului *lucra* din limba română.

Astfel se poate afirma că este vorba de un schimb de prefixe, menționându-se, în aceeași ordine de idei, și faptul că nu sunt așa teme precum *păduri*, și *vălui*. Un moment interesant este că în dicționarul morfologic s-au găsit doar verbe care reprezintă derivate semianalizabile, chiar dacă, în afară de *despăduri*, mai este și substantivul *despădurire*, care nu este găsit ca fiind semianalizabil. În cazul derivatelor resursa www.dexonline.ro indică cuvântul de la care a fost derivat, prin marcarea cu simbolul *V.* care ar însemna „vezi cuvântul”.

4. Proiectarea derivării cu prefixe pe ulterioara derivare cu sufixe

Proiectarea derivatelor reprezintă o metodă de formare a cuvintelor prefixate de la derivate sufixate de la aceeași rădăcină. Potrivit cercetătorilor spanioli (Santana et al., 2004), verbul spaniol *amortizar* poate fi derivat cu prefixul *des-* în *desamortizar*. Totodată *amortizar* poate fi derivat cu sufixele *-cion* și *-able*. Astfel, derivatul cu prefixul *des-* poate deriva cu sufixele *-cion* și *-able*. Apare ipoteza potrivit căreia derivatele cu prefixe pot moșteni/proiecta derivatele cu sufixe ale temei de la care a fost realizată derivarea cu prefixe. Deci, mecanismul de generare a derivatelor ar presupune următoarele: dacă *R* este rădăcina cuvintelor, S_i sufixele posibile pentru rădăcina *R*, adică $R \rightarrow RS_i$, și *P* prefixul corespunzător rădăcinii *R*, adică $R \rightarrow PR$, atunci există sufixe S_i pentru care $R \rightarrow PRS_i$

Prin urmare, dacă luăm drept exemplu tema *capitula*, atestăm conform www.dexonline.ro următoarele derivate cu sufixe: *capitulant*, *capitulantă*, *capitulard*, *capitulare*, *capitulație*. Totodată, pentru derivatul cu prefixul *re-*, *recapitula* are următoarele derivate cu sufixe: *recapitulare*, *recapitulație*. Astfel, sunt două derivate care realizează o proiectare a derivării cu sufixe. Încercând proiectarea derivatului *capitulant* \rightarrow *recapitulant*, se obține un cuvânt inexistent, chiar și pentru resursele electronice ale motorului de căutare www.google.com.

Pornind de la lexiconul existent cu structura morfemică (Constantinescu, 2008), cu ajutorul unor programe, s-au extras acele grupuri de derivate la care se atestă proiectarea, înregistrându-se 363 de teme de la care este posibilă proiectarea derivatelor.

Marea majoritate a acestora constă doar dintr-un derivat sufixat și unul prefixat. Totuși, se poate afirma că metoda este utilă în generarea derivatelor pentru limba română. Un exemplu de astfel de derivare este ilustrat în Figura 2.

5. Studiarea diminutivelor și a augmentativelor, derivarea prin schimbarea genului

Potrivit cercetătorilor sârbi (Duško&Krstev, 2005), cu ajutorul cuvintelor derivate din limba sârbă, au fost generate noi leme cu sensuri previzibile. Prin sensuri previzibile se înțelege derivarea prin amplificarea sensului, adică generarea diminutivelor (*profesorčić*) și a augmentativelor (*profesorčina*), schimbarea genului cuvântului (pentru cel masculin *professor* → cel feminin *professorka*),

Procesarea diminutivelor

Sufixe diminutive ar fi următoarele: *-aș* (copilaș), *-uc* (sătuc), *-el* (bătrânel), *-iță* (fetiță), *-uță* (caruță), *-ică* (floricică), *-uleț* (ursuleț), *-iș* (podîș), *-uț* (căluț).

Tabel 2: Numărul de derivate diminutive pentru sufixe concrete

Sufixul	Numărul de derivate
<i>-aș</i>	327
<i>-iță</i>	249
<i>-el</i>	221
<i>-uță</i>	208
<i>-ică</i>	139
<i>-uleț</i>	104
<i>-iș</i>	101
<i>-uț</i>	88
<i>-uc</i>	7

Tabel 3: Numărul de derivate augmentative pentru sufixe concrete

Sufixul	Numărul de derivate
<i>-andru</i>	4
<i>-an</i>	74
<i>-oi</i>	74
<i>-oaie</i>	26

Potrivit aceleași surse (Constantinescu, 2008), au fost stabilite numărul de derivate cu fiecare sufix diminutival (Tabel 2.). Dincolo de aceasta, toate aceste sufixe înregistrează unele alternanțe vocalice și/sau consonantice în cazul procesului de derivare, de exemplu, *sat* - *sătuc*, *car* - *căruță*, *cal* - *căluț* unde alternanța este *a->ă*; *fată* - *fetiță*: *a ->e*; *floare* - *floricică*: *oa ->o*; *frate*- *frățior*: *at -> ăț* etc.

Cea mai numeroasă clasă formată în urma procesului de derivare cu aceste sufixe este clasa substantivelor, menționându-se aici următoarele sufixe *-uc*, *-el*, *-aș*, *-iș*, *-uț*, *-uță*, *-iță*, *-uleț* și *-ică*. Mai puțin numeroase sunt adjectivele, în cazul sufixelor: *-uc*, *-el*, *-iș*, *-uț* și *-ică*. Un caz foarte rar s-a atestat cu sufixul *-iș*, care a format de la verbul *zbură* adverbul *zburîș*.

Procesarea augmentativelor

Sufixe augmentative sunt următoarele: *-andru* (copilandru); *-an* (băietan); *-oi/oaie* (căsoi, căsoaie). Potrivit lexiconului (Constantinescu, 2008), au fost stabilite numărul de derivate cu fiecare sufix augmentativ (Tabel 3.). Toate aceste sufixe pot fi atașate substantivelor pentru a obține substantive. Sufixele *-an*, *-oi*, *-oaie* pot fi atașate

adjectivelor pentru a obține cuvinte augmentative. De menționat că în procesul de derivare cu sufixele respective se atestă unele alternanțe vocalice și/sau consonantice, cum ar fi de exemplu, *casă* - *căsoi/căsoaie*, unde alternanța este *a->ă*, *băiet* - *băiețandru* *t->ț*, etc.

Schimbarea genului cuvântului

În cazul limbii române, schimbarea genului cuvântului poate fi realizată prin trecerea de la unele sufixe la altele corespunzătoare, de exemplu, *-tor* ↔ *-toare*, *-esc* ↔ *-ească*, etc. Astfel, se observă că la schimbarea genului se apelează la procesul de sufixare, și nu la cel de prefixare.

Lexiconul menționat anterior conține derivate care sunt sufixate cu *-tor*, cu *-toare* și atât cu *-tor* cât și cu *-toare*. Astfel, potrivit informației din Figura 3, există 148 de cuvinte (substantive și/sau adjective) de forma $\omega'=\omega tor$, care ar putea deriva în cuvinte de forma $\omega''=\omega toare$. Similar, sunt 42 de cuvinte (substantive și/sau adjective) de forma $\beta'=\beta toare$ care e posibil să deriveze în $\beta''=\beta toare$. Totodată, aceste 190 de derivate generate automat urmează să fie validate. Pentru început cuvintele au fost verificate la prezența lor în RRTLN. Din toate cuvintele generate 122 au fost prezente pe când restul cuvintelor au fost verificate în documentele electronice din Internet, din 68 de derivate, fiind validate 49. Prin urmare 95% din cuvintele generate au fost valide.

Aceeași situație este și cu perechea de sufixe *-esc* și *-ească*. Potrivit aceluiași lexicon (Constantinescu, 2008), se conțin 274 de derivate cu sufixul *-esc*, și 249 cu *-ească*. De menționat că sunt 229 de derivate care sunt sufixate atât cu *-esc* cât și cu *-ească*. Devine firească ipoteza potrivit căreia cuvintele (substantive și/sau adjective) de forma $\omega'=\omega esc$, pot deriva în cuvinte de forma $\omega''=\omega ească$. Similar, cuvintele (substantive și/sau adjective) de forma $\beta'=\beta ească$ pot deriva în cuvinte de forma $\beta''=\beta esc$. Generând automat acele derivate care lipsesc în cazul genului și verificând automat prin existența cuvintelor în documentele electronice, s-a constatat că cu ajutorul RRTLN au fost validate 43 de cuvinte din 65 generate. Din celelalte 22 de derivate s-a reușit de a valida încă 12 derivate cu ajutorul unei aplicații Web, bazate pe posibilitățile motorului de căutare Google. În general s-au obținut 84% de cuvinte valide.



Figura 3: Numărul derivatelor cu sufixele *-tor* și *-toare*

6. Unele aspecte în generarea derivatelor cu *in-/im-*

Există mai multe clase de derivate cu prefixele *in-/im-*. În cele ce urmează se va descrie cazul în care se formează derivate cu aceste prefixe cu sens negativ. În limba română sunt mai multe prefixe care oferă derivatului sensul de negație, precum ar fi: *a-*, *i-*, *ne-* și *im-/in-*.

Derivatele cu *im-/in-*, de obicei, sunt adjective, rareori substantive sau verbe. Cele mai numeroase derivate cu prefixul *in-/im-* sunt raportate la adjective formate cu sufixul -

bil, de exemplu, *incurabil*, *inestimabil*, etc. Așadar, fiind date adjectivele $\omega' = \omega bil$ se formează derivatele $\omega'' = \beta \omega bil$, unde $\beta \in \{in-, im-\}$.

O altă clasă bine reprezentată este cea a formațiilor raportabile la adjective derivate cu sufixul *-ent* și *-ant*: *inaderent*, *incoerent*, *independent*, etc. (Iorgu, 1970). La fel, fiind date adjectivele $\omega' = \omega \gamma$, se formează derivatele de forma $\omega'' = \beta \omega \gamma$, unde $\beta \in \{in-, im-\}$ și $\gamma \in \{-ent, -ant\}$. În ambele cazuri alegerea lui β depinde de prima literă a adjectivului ω , și anume în cazul în care această literă este b sau p atunci $\beta = im-$, în alte cazuri e $\beta = in-$.

Alte clase de derivate cu *in-/im-* sunt neînsemnate. De menționat că există derivate în *im-/in-* de la teme deja derivate cu prefixe.

Potrivit resursei www.dexonline.ro, în limba română sunt circa 4946 de cuvinte care încep cu combinația de litere *in* și 1249 - cu combinația de litere *im*. Din acestea care se termina cu *ant* sunt corespunzător 38 și 13, cu *ent* 61 și 12, și respectiv *bil* 220 și 43, în total fiind 387 de cuvinte. Așa precum unele cuvinte, cum ar fi cele ce se termina cu *bil*, pot aparține mai multor părți de vorbire, s-a reușit filtrarea și obținerea unui număr de doar 293 de adjective. Mai mult decât atât, adjectivele cu sufixul *-bil*, mai formează și derivate cu prefixul *ne-*, care, de altfel, oferă același sens negativ. În această ordine de idei, este utilă verificarea derivatelor formate cu ajutorul motorului de căutare, de exemplu www.google.com.

Cazul derivării cu prefixul *in-/im-* de la cuvinte care se termină în *-ant* sau *-ent* este și mai problematic, prin faptul că sunt multe cuvinte care nu vor forma derivate în *in-/im-* pe motiv că sunt și substantive, nu doar adjective.

Lexiconul de derivate (Constantinescu, 2008) conține doar un singur derivat cu prefixele *im-/in-*, și anume *impermeabilizare*. Astfel, este posibil de a genera derivate cu prefixele corespunzătoare. În mod automat, s-a constatat că lexiconul (Constantinescu, 2008) conține 62 de derivate cu *-bil*, 1 derivat cu *-ent* și 37 cu *-ant*.

Examinând cuvintele din lexicon și concatenându-le prefixul respectiv celor care se încadrează în categoriile stabilite, fără alternanțe vocalice sau consonantice, s-a construit un algoritm de derivare cu prefixul *im-/in-*, în baza căruia a fost elaborat un modul care generează cuvinte noi. Drept rezultat s-au obținut 100 de derivate cu prefixele *in-/im-*.

Verificând inițial cuvintele generate la prezența lor în RRTLN, s-a constatat prezența a 7 cuvinte. Celelalte 93 de cuvinte au fost verificate în documentele electronice din Internet cu ajutorul motorului de căutare www.google.com. Ca rezultat, doar 14 derivate din toate cele generate au fost regăsite în documentele electronice. Patru cuvinte s-au regăsit o singură dată, și anume: *inofertant*, *imprelucrabil*, *inrezolvabil*, *intrasabil*. Datele pentru restul derivatelor sunt prezentate în Tabel 4.

Tabel 4: Numărul de apariții a derivatelor generate în documentele electronice Internet

Derivatele	Numărul de apariții
<i>Invindecabil</i>	7
<i>indefrișabil</i>	8
<i>Indifuzabil</i>	8
<i>Infiltrabil</i>	8

Derivatele	Numărul de apariții
<i>indeșirabil</i>	56
<i>Insubstituibil</i>	77
<i>Imprecizabil</i>	94
<i>Injustificabil</i>	181
<i>Injucabil</i>	353
<i>Incadabil</i>	4710

În cazul derivatelor generate cu prefixul *in/im* cu un număr mic de apariții printre documentele electronice, este clară poziția potrivit căreia nu a fost aplicat prefixul corect cu sens negativ, aceasta fiind situația derivatelor *inrezolvabil* (1) – *nerezolvabil* (1280), *insubstituibil* (77) – *nesubstituibil* (222), *injucabil* (353) - *nejucabil* (3050). Totodată, mai sunt derivate cu un număr mare de apariții și care au variante cu alt prefix și cu un număr mai mic de apariție a aceluși prefix, aici fiind depistate următoarele derivate: *inabordabil* (2810) - *neabordabil* (699), *inacceptabil* (67900) - *neacceptabil* (7140), *incalculabil* (24000) - *necalculabil* (469)

7. Concluzii

Generarea derivatelor nu este o problemă deloc banală, deoarece procesul de derivare a cuvintelor nu prezintă un mecanism regulat. Soluția potrivit căreia se pot stoca toate derivatele într-un dicționar este una rezonabilă, or, aceste derivate înregistrate corect oricum nu vor acoperi întreaga diversitate a limbii, care presupune un mecanism în continuă dezvoltare. Pe de altă parte, abordarea potrivit căreia se poate genera derivate conform unor reguli de constrângere a grupurilor de derivate, este un mecanism de supragenerare, în care etapa de validare exclude foarte multe cuvinte formate greșit. Doar stabilirea unor reguli bine puse la punct ar putea ridica nivelul de generare a cuvintelor corecte.

Referințe bibliografice

- Constantinescu, S. (2008). *Dicționar de cuvinte derivate*, Editura Herra, București.
- Duško, V., Krstev, C. (2005). Derivational Morphology in a E-Dictionary of Serbian, In Zygmunt Vetulani (ed.), *Proceedings of the 2nd Language & Technology Conference*, Poznan, Poland, pp. 139-143.
- Iordan, I. (1970). *Limba română contemporană*. București, pp. 66-99.
- Petic, M. (2010). Automatic derivational morphology contribution to Romanian lexical acquisition. *Special issue: Natural Language Processing and its Application. Research in Computing Science*, Mexico, vol. 46, pp. 67-78.
- Santana, O., Perez J., Carreras, F., Rodrigues, G. (2004). Suffixal and Prefixal Morphological Relationships of Spanish. *Lecture Notes in Artificial Intelligence*, Ed. Springer-Verlag, pp. 407-418.

DEZVOLTAREA UNUI PARSER DE ROLURI SEMANTICE PENTRU LIMBA ROMÂNĂ

DIANA TRANDABĂȚ¹, DAN CRISTEA^{1,2}

¹*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

²*Universitatea „Al. I. Cuza”, Facultatea de Informatică, Iași*

{*dtrandabat, dcristea*}@*info.uaic.ro*

Rezumat

Parsarea semantică, prin identificarea și clasificarea entităților semantice în context, precum și a relațiilor dintre ele, are un mare potențial pentru aplicații cum ar fi rezumare de text, sisteme de întrebare-răspuns, sau traducere automată. Astfel, prin dezvoltarea unui sistem care să adnoteze automat roluri semantice pentru limba română, lucrarea de față reprezintă un pas intermediar important pentru înțelegerea automată a limbajului natural. Pentru dezvoltarea parserului de roluri semantice pentru limba română, a fost întâi necesară dezvoltarea unui corpus de antrenare, adnotat la roluri semantice. Această resursă adnotată de roluri semantice pentru limba română a fost creată prin utilizarea unei metode de transfer automat al rolurilor semantice, plecându-se de la o resursă dezvoltată pentru limba engleză (FrameNet). Ulterior, folosind o platformă pentru dezvoltarea de sisteme de etichetare supervizată a rolurilor semantice (PASRL), au fost antrenați mai mulți algoritmi de învățare folosindu-se corpusul dezvoltat, iar modelul obținut cu cel mai bun algoritm a fost salvat. Sunt discutate rezultatele aplicării acestui model pe texte neadnotate din limba română.

1. Introducere

Semantica cadrelor (*Frame Semantics*) (Fillmore, 1982) este o teorie lingvistică care descrie structura conceptuală ce stă la baza înțelesului lingvistic. Cadrul semantic reprezintă o structură scriptică de inferențe, legată prin convenții lingvistice de înțelesul unităților lexicale. Fiecare cadru identifică un set de constituenți, roluri semantice (*frame elements*) care îl definesc și o serie de unități lexicale (cuvinte) ce participă la actualizarea sa. Conceptul de *unitate lexicală* (sau cuvânt predicțional) este central pentru resursele de cadre semantice. Unitatea lexicală este acel cuvânt sau sens al unui cuvânt polisemantic pentru care se definesc proprietățile combinatorice, definit printr-o leamă, o parte de vorbire și un cadru. Descrierea în termeni de semantică a cadrelor a unei unități lexicale identifică cadrele care formează un înțeles dat și specifică modul în care rolurile semantice sunt realizate în interiorul unor structuri dominate de cuvântul țintă.

Sistemul prezentat este dedicat prelucrării limbajului natural plecând de la text neadnotat, cu scopul de a determina informațiile de natură semantică pentru verbele și substantivele predicționale. Cuvintele predicționale au o anumită valență, și anume capacitate combinatorie, însușire de a deschide anumite poziții libere care sunt ocupate de termenii învecinați în propoziție. Fiind strâns legată de semantica verbului, valența determină numărul și caracteristica funcțional-semantică și gramaticală a elementelor cerute de verb (constituenți esențiali și facultativi). Cel mai activ și mai important sub

aspectul valenței este verbul, însă pot exista și substantive sau adjective ce au trăsătura de predicționalitate care realizează nucleul semantic al propoziției. Sistemul de adnotare automată a rolurilor semantice pentru limba română se bazează pe o resursă de cadre semantice creată pentru limba română prin transferul rolurilor din limba engleză (Trandabăț, 2007).

Lucrarea este structurată în 4 secțiuni. După o scurtă introducere în domeniul rolurilor semantice, secțiunea 3 începe cu descrierea modului de creare a unui corpus de antrenare adnotat la roluri semantice pentru limba română. Folosirea acestui corpus pentru dezvoltarea unui sistem de adnotare automată a rolurilor semantice este prezentată în secțiunea 3.2., iar secțiunea 3.3. prezintă succint modulele de pre-procesare necesare pentru folosirea parserului de roluri semantice pentru limba română. În secțiunea 4 sunt discutate direcții viitoare de dezvoltare și utilizare a sistemului de etichetare a rolurilor semantice pentru limba română.

2. Rolurile semantice și adnotarea lor automată

În procesul comunicării, cuvintelor predicționale sunt completate cu constituenți adiționali. Gradul de necesitate al constituenților pentru plenitudinea semantică a enunțului alcătuit de predicat este însă diferit. De exemplu, verbul *a depinde* cere cu precădere un *actant* și un *obiect* (*cineva depinde de ceva*). Ceilalți determinați care pot apărea pe lângă verbul *a depinde* sunt actualizați explicit numai atunci când sunt ceruți de situația de comunicare. Astfel, unele dintre pozițiile sintactice descrise de verb pot rămâne libere, fără ca enunțul să aibă de suferit și, dimpotrivă, suprimarea sau neexprimarea altora poate avea drept efect generarea unor structuri incorecte, incomplete sub aspect structural și semantic. Din acest punct de vedere se disting două tipuri de valențe (și deci tipuri de roluri semantice): valențe esențiale, numite și *argumente* și valențe facultative, numite și *adjuncți*. Un rol semantic esențial (numit și rol nucleu) este un rol care instanțiază o componentă necesară din punct de vedere conceptual pentru definirea și diferențierea cadrului. *Cumpărător*, *Vânzător*, *Bani* și *Bunuri* sunt exemple de roluri nucleu pentru cadrul semantic *Comerț*. Rolurile semantice care introduc evenimente independente sau distincte, adiționale evenimentului principal, sunt facultative, și se mai numesc și roluri periferice. Rolurile periferice marchează dimensiunea temporală, spațială, modală etc., dar nu caracterizează în mod unic cadrul și pot actualiza orice cadru semantic evenimential. *Mod*, *Mijloace*, *Scop*, *Rată* și *Unitate* sunt exemple de roluri periferice pentru cadru semantic *Comerț*.

Etichetarea rolurilor semantice este o componentă importantă a înțelegerii limbii, și a fost considerată de mai multe sisteme computaționale. Sistemele tradiționale de parsare și înțelegere, inclusiv implementări bazate pe gramatici de unificare, se bazează pe gramatici dezvoltate manual, care trebuie să anticipeze fiecare mod în care rolurile semantice ar putea fi realizate sintactic. Scrierea acestor gramatici este consumatoare de timp, și de obicei, astfel de sisteme au o performanță limitată. Metodele bazate pe învățare promit o generalizare dincolo de numărul relativ mic de instanțe sau roluri considerate. Astfel, diverse strategii de învățare au fost folosite pentru adnotarea automată a rolurilor semantice: estimarea probabilităților (Gildea and Jurafsky, 2002),

arbori de decizie (Surdeanu et al., 2003), mașini cu suport vectorial (Pradhan et al., 2005) și învățarea bazată pe memorie (Morante et al., 2008).

Un dezavantaj important al celor trei sisteme prezentate este că acestea nu tratează predicatul nominal, fiind construite doar pentru predicatul verbal¹. Mai mult, acestea iau în considerare doar un singur predicat pentru fiecare propoziție, chiar dacă nu acesta este întotdeauna cazul. De exemplu, în propoziția:

Acordarea premiului Nobel Președintelui Obama a fost dezbătută pe larg.

avem două cuvinte predicative, acordarea, având ca *Temă* premiului Nobel, și dezbătute, având două roluri, o *Temă* acordarea Premiului Nobel Președintelui Obama și un adjunct *Modal* reprezentat de grupul prepozițional pe larg.

Sistemul prezentat în lucrarea de față tratează atât verbele, cât și substantivele predicative pentru limba română.

3. Dezvoltarea parserului de roluri semantice pentru limba română

3.1 Corpusul de antrenare

Pentru dezvoltarea resurselor de cadre semantice pentru limbile spaniolă, germană și japoneză, s-a plecat de la un corpus specific fiecăreia dintre aceste limbi, adnotat manual la roluri semantice. În această secțiune descriem o metodă de constituire a unui corpus românesc de cadre semantice prin import din limba engleză. Premisa programului de importare automată a rolurilor semantice din limba engleză pentru limba română se bazează pe proprietatea cadrelor semantice de a exprima concepte la nivelul structurii de adâncime, valabile pentru toate limbile, actualizarea sintactică având loc ulterior la nivelul structurii de suprafață, diferit pentru fiecare limbă în funcție de constrângerile sintactice și morfologice. Programul de transfer automat (Trandabăț, 2007) are la bază corelarea rolurilor semantice exprimate în limba engleză cu traducerea pentru limba română a cuvintelor ce realizează rolurile respective. Ulterior, aceste adnotări automate au fost validate manual. Această abordare este similară celei din (Barbu Mititelu and Ion, 2005) folosită pentru transferul relațiilor de dependență verbală dintr-un corpus aliniat englez-român.

Astfel, folosind fișierele XML ale propozițiilor engleze adnotate la roluri semantice, se creează automat un set de fișiere XML ce conțin un corpus de propoziții adnotate la nivelul rolurilor semantice pentru limba română, din care urmează să se extragă cadrele semantice. Programul de adnotare folosește ca fișiere de intrare: (i) fișierele XML pentru unitățile lexicale engleze, care conțin propoziții adnotate, și (ii) fișierele cu alinierea propozițiilor engleze - românești.

Pentru realizarea corpusului românesc au fost alese aleatoriu 110 de propoziții din resursa FrameNet pentru limba engleză. Pe lângă aceste propoziții, a fost ales cadrul semantic *Event* cu toate cadrele semantice corelate lui, totalizând alte 984 de propoziții

¹ Începând cu competiția ConLL2008, predicatul nominal a fost introdus în sistemele de identificare a rolurilor semantice, împreună cu predicatul verbal.

adnotate. În plus, încă 400 de propoziții, obținute din competiția de aliniere la nivel de cuvânt pentru limbile engleză-română din cadrul ACL 2003 și 2005², au fost adnotate cu un parser de roluri semantice pentru limba engleză, și ulterior validate.

După selectarea propozițiilor s-a realizat traducerea manuală a enunțurilor (exceptând cele 400 de propoziții deja traduse) încercându-se menținerea pe cât posibil a părții de vorbire din limba engleză, măcar pentru unitatea lexicală care determină cadrul semantic, pentru a mări precizia alinierii.

Alinierea propozițiilor englezești cu cele românești a fost realizată folosindu-se aliniatorul dezvoltat de Institutul de Cercetări în Inteligență Artificială al Academiei Române (Tufiş et al., 2005).

Au fost identificate 9 tipuri de import al adnotării, descrise detaliat în (Trandabă, 2010). Validarea rezultatelor s-a bazat pe detectarea cazurilor când importul a eșuat, încercând să descopere dacă problemele s-au datorat traducerii sau particularităților semantice și sintactice ale limbii române. Au fost găsite doar puține erori de traducere, și chiar și în acele cazuri, înțelesul fusese păstrat și rolurile semantice fiind corect atribuite.

Majoritatea propozițiilor constituie un corpus corect din punct de vedere al adnotării. Cazurile speciale de neconcordanțe țin fie de inconsecvențe de adnotare, fie de diferențe culturale sau lingvistice dintre cele două limbi. Problemele de import au fost prezentate în (Trandabă & Husarciuc, 2008) și pot fi grupate în următoarele clase:

- Cazuri în care există mai multe adnotări posibile în engleză – același constituent poate avea roluri tematice diferite. În limba română pot exista ambele roluri sau doar unul. Programul de import recunoaște un singur rol pentru fiecare constituent, și anume primul care a fost adnotat. Aceste cazuri sunt marcate pentru validare manuală. De exemplu, importul automat pentru propoziția din limba engleză:

Traditional methods require that [the animal]_{Protagonist} [bleed]_{Cause} to [death]_{TARGET} [after having its throat cut]_{Time/Cause}

este:

Metodele tradiționale cer ca [animalul]_{Protagonist} [să sângereze]_{Cause} până la [moarte]_{TARGET} [după ce i s-a tăiat gâtul]_{Time}.

- Cazuri în care limba română are roluri semantice ce nu apar în limba engleză (de exemplu argumente externe) ca în:

[Quit]_{TARGET} [smoking]_{Process}.

[Lăsați]_{TARGET} –[vă]_{Protagonist} [de fumat].

unde rolul *Protagonist* nu există în propoziția engleză, dar el apare în traducerea în limba română și ar trebui adnotat, deși programul de import nu poate transfera nimic.

- Cazuri în care un rol din limba engleză nu este exprimat în limba română deoarece nu este exprimat explicit în structura sintactică (de obicei este vorba de subiect care este

² Propozițiile au fost descărcate de pe pagina Radei Mihalcea: <http://www.cse.unt.edu/~rada/downloads.html#romanian>

obligatoriu în limba engleză, dar poate lipsi fără a invalida enunțul în limba română). Un exemplu în care rolul *Manner* este inclus în cuvântul țintă este:

[Blood]_{Undergoer} [had congealed]_{TARGET} [thickly]_{Manner} [on the end of the smashed fibula]_{Place}.

[Sângele]_{Undergoer} [se îngroșă]_{TARGET} [spre capătul fibulei zdrobite]_{Place}.

- Diferențe privind modul de formulare a enunțurilor în cele două limbi. Cele mai multe exemple de acest fel găsite au fost datorate folosirii în limba engleză a verbelor copulative, modale, sau a verbelor support, traduse în limba română fără folosirea unui astfel de verb. Astfel, cadrele semantice din cele două limbi sunt diferite, iar rolurile din limba engleză, deși importate corect, pot fi inexistente în cadrul semantic al verbului din limba română.

And I was surprised at how easily [my eyes]_{Entity} [became]_{TARGET} [accustomed to seeing]_{Final_state} in the light of the head torch.

Și am rămas surprins cât de ușor [s- au obișnuit]_{TARGET} [ochii mei]_{Entity} [cu vederea]_{Final_state} la lumina făcliei principale.

Adnotările automate ale rolurilor semantice pe corpusul românesc au fost verificate pentru a se extrage situațiile de neconcordanță. O dezvoltare ulterioară este realizarea, plecând de la cazurile de diferențe lingvistice, de reguli automate ce vor fi implementate pentru îmbunătățirea rezultatelor programului de transfer automat al rolurilor semantice.

În plus față de adnotarea semantică, corpusul a fost îmbogățit cu adnotare la partea de vorbire, folosindu-se serviciul web RACAI (Tufiș et al., 2008), și adnotare a dependențelor sintactice, folosindu-se un parser creat prin antrenarea MALTParser (Nivre, 2003) pe un corpus de fraze românești adnotate manual cu dependențe FDG.

3.2 PASRL

Folosind resursa descrisă în secțiunea precedentă, a fost dezvoltat un sistem de creare a unui modul de adnotare automată a rolurilor semantice, numit PASRL (Platform for Adjustable Semantic Role Labeling), descris în (Trandabăț, 2010). Similar cu arhitectura generală a sistemelor de etichetare a rolurilor semantice (Marquez et al., 2008), PASRL este compus din două sub-sisteme principale: un modul de predicție a predicatului și un modul de predicție a argumentelor. Modulul de predicție a predicatului are la rândul lui două configurații posibile, corespunzând identificării predicatelor și a sensurilor acestor predicate succesiv sau simultan (figura 1).

Identificarea predicatelor – acest modul primește ca intrare propoziția analizată sintactic (cu informații despre părțile de vorbire și dependențele sintactice) și decide care dintre verbele și substantivele din propoziție sunt predicabile, și deci pentru care trebuie căutate roluri semantice;

Identificarea sensului predicatelor – după identificarea predicatelor dintr-o propoziție, trebuie stabilit sensul fiecărui predicat (dintre sensurile din corpusul de antrenare), deoarece sensuri diferite pot cere o structură diferită de roluri semantice;

Identificarea rolurilor semantice – acest modul identifică, pentru fiecare predicat, ce rol are fiecare dependent al predicatului.

DEZVOLTAREA UNUI PARSER DE ROLURI SEMANTICE PENTRU LIMBA ROMÂNĂ

Pentru fiecare sub-problemă (nivel din figura 1), modulele au trei variante, legate de dimensiunea corpusului de antrenare: folosind toate datele de test, folosind numai datele care au drept cuvânt predicational un substantiv, respectiv verb (NP / VP), sau folosind numai datele corespunzând unui anumit tip de rol. Pentru fiecare modul din figura 1, sunt antrenați un set de algoritmi de clasificare din Weka (Witten and Frank, 2005). După rularea tuturor clasificatorilor pentru toate modulele, performanța lor este comparată, iar calea care obține cea mai mare performanță este considerată cea mai bună configurație. Un exemplu de astfel de configurație poate fi: Rularea modelului creat pentru identificarea predicatelor folosind tot corpusul de antrenare și drept algoritm de clasificare arboreii de decizie, urmat de modelul creat pentru identificarea sensului predicatului separat pe verbe/substantive (NP / VP) create tot folosind arbori de decizie, și de modelul de identificare a argumentelor pentru fiecare tip de argument în parte, creat folosind algoritmul Naive Bayes. Modelele pentru cea mai bună configurație sunt salvate iar configurația este folosită ulterior pentru a adnota texte noi cu roluri semantice.

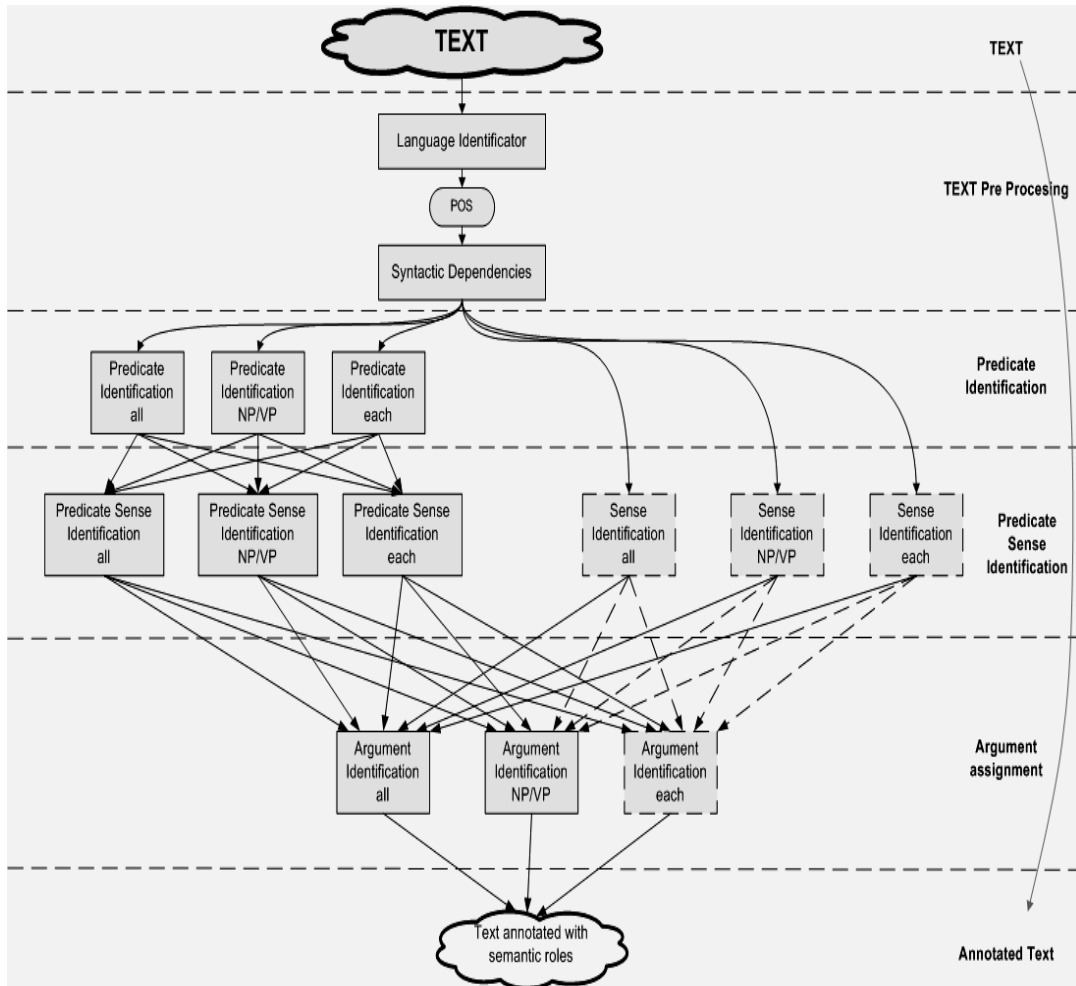


Figura 1: Arhitectura sistemului PASRL

Pentru stabilirea instanțelor, pentru fiecare predicat sunt considerați dependenții lui sintactici și contextul său imediat (5 cuvinte înainte și după predicat). Trăsăturile folosite în modulului de identificare a predicatului și pentru modulul de identificare a sensului predicatului sunt extrase din resursa creată pentru limba română:

1. o trăsătură binară care stabilește dacă cuvântul predicțional se află sau nu în lista de predicate adnotate în resursa pentru limba română³;
2. numărul de dependenți sintactici (dacă un cuvânt nu are dependenți, atunci nu contează dacă el este sau nu predicțional, pentru că nu este considerat predicat pentru nici un constituent din propoziție⁴);
3. numărul de dependenți sintactici situați înainte/după predicat (de ex., mai mulți dependenți înainte de predicat indică o topică anormală sau un verb la diateza pasivă, deci rolurile semantice pot fi inversate);
4. părțile de vorbire și grupurile din care fac parte n cuvinte înainte, respectiv după cuvântul predicat, în ordinea de suprafață a propoziției, unde n are valoarea 5, stabilită empiric, dar poate fi modificată dacă este dată ca parametru;
5. relațiile de dependențe sintactice a celor n cuvinte;
6. trăsătură binară care stabilește dacă părintele dependentului sintactic este predicatul sau nu;
7. trăsătură binară care stabilește dacă părintele dependentului sintactic este unul dintre cele n cuvinte;
8. trăsătură binară care stabilește dacă predicatul are dependenți care nu se găsesc între cele n cuvinte.

Pentru modulul de identificare și clasificare a rolurilor semantice, instanțele sunt calculate considerându-se toate cuvintele din propoziție și următoarele trăsături:

1. partea de vorbire a cuvântului;
2. dependența sintactică față de predicat;
3. grupul sintactic din care face parte;
4. partea de vorbire a regentului cuvântului;
5. distanța de la cuvânt până la predicat, înainte sau după (în număr de cuvinte);
6. trăsătură binară stabilind dacă cuvântul investigat este punctuație sau nu;
7. trăsătură binară stabilind dacă cuvântul investigat este fiu al predicatului sau nu;
8. trăsătură binară stabilind dacă regentul cuvântului investigat este „a fi”, un verb copulativ sau modal;
9. hipernimul din WordNetul românesc (Tufiş et al., 2006), pentru substantive. Hipernimele sunt clasificate în prima clasă găsită, mergând bottom-up în ierarhia hipernimelor, dintre: spațial, temporal, cantitate, obiect, persoană, entitate.

PASRL a fost rulat pentru limba română, și cea mai bună configurație a fost: identificarea predicatelor cu arbori de decizie folosind tot corpusul, identificarea sensurilor predicatelor folosind algoritmul SimpleCart antrenat separat pentru predicate verbale, respective nominale, urmat de identificarea argumentelor folosind arbori de

³ Resursa de roluri pentru limba română permite crearea unei liste a predicatelor (și a sensurilor acestora) pentru care există adnotări. Această listă este momentan restrânsă numeric, dar o dată cu îmbogățirea resursei românești, ea va fi actualizată, această trăsătură devenind mai importantă pentru determinarea acelor cuvinte care pot avea roluri.

⁴ În această versiune a parserului, nu sunt adnotate rolurile implicite sau cele neexprimate.

decizie. Evaluarea performanțelor folosind 10-fold cross-validation pe corpusul de antrenament indică o precizie de 74% și un recall de 64%.

3.3 Ajustarea parserului pentru text neadnotat

Pentru a adnota un text nou în limba română cu roluri semantice, utilizând modele salvate de PASRL, informațiile sintactice sunt adăugate folosind serviciile web RACAI, iar relațiile de dependență folosind un model antrenat pentru limba română cu parserul MALT. Evaluările preliminare indică o F-measure de 54%, ceea ce indică un început promițător, având în vedere dimensiunea redusă a corpusului pe care s-a efectuat învățarea (1500 de propoziții).

4. Concluzii și direcții viitoare

Intuiția conform căreia analiza semantică poate reprezenta o contribuție pozitivă majoră pentru aplicațiile de prelucrarea limbajului natural a motivat dezvoltarea unui număr de resurse de semantică lexicală care cuprind descrierea cadrelor semantice ale fiecărui verb dintr-o limbă. Cele mai importante astfel de proiecte sunt PropBank (Palmer et al., 2005) și FrameNet (Baker et al., 1998). Valoarea acestor resurse este direct dependentă de cantitatea de informații pe care le conțin, lucru de asemenea dependent de nivelul efortului implicat în dezvoltarea lor. Această lucrare prezintă un corpus românesc de roluri semantice realizat printr-o metodă de transfer și un sistem de adnotare a rolurilor semantice pentru limba română. Una din principalele direcții urmărite în viitor este îmbunătățirea resursei de roluri semantice (și implicit și a sistemului de etichetare automată) prin transferarea altor propoziții din limba engleză. De asemenea, metoda de import poate fi folosită și pentru alte domenii semantice. În această direcție, vom investiga transferarea adnotărilor referitoare la referințele anaforice, urmând testele din (Postolache et al., 2006).

Sistemul de etichetare a rolurilor semantice pentru limba română va fi în curând disponibil sub formă de serviciu web, pentru a putea fi implicat în dezvoltarea sistemelor aplicative de prelucrare a limbajului natural pentru limba română, cum sunt sistemele de rezumare automată (Cristea et al., 2005) sau de întrebare-răspuns (Iftene et al., 2009), dar și pentru a fi folosite în proiecte de integrare a limbii române în rețele europene (Cristea și Pistol, 2009).

Mulțumiri. Autorii mulțumesc colectivului Institutului de Cercetări în Inteligența Artificială al Academiei Române pentru accesul public la serviciile web. De asemenea, mulțumim colegilor Augusto Perez și Alex Moruz pentru adnotarea propozițiilor FDG, respectiv antrenarea parserului de dependențe sintactice pentru limba română. Cercetarea prezentată în această lucrare a fost finanțată de către Programul Operațional Sectorial Dezvoltarea Resurselor Umane prin proiectul „Dezvoltarea capacității de inovare și creșterea impactului cercetării prin programe post-doctorale POSDRU/89/1.5/S/49944.

Referințe bibliografice

- Baker, C. F., Fillmore, C. J., and John B., L. (1998). *The Berkeley FrameNet project*. In Proceedings of the COLING-ACL, Montreal, Canada.
- Barbu Mititelu, V. and Ion, R. (2005). *Automatic Import of Verbal Syntactic Relations Using Parallel Corpora*. In Proc. of the International Conference Recent Advances in Natural Language Processing, pages 329-333, Borovets, Bulgaria.
- Cristea, D., Postolache, O., and Pistol, I. (2005). *Summarisation through discourse structure*. In Alexander Gelbukh, editor, Computational Linguistics and Intelligent Text Processing, 6th International Conference CICLing 2005, volume 3406 of LNCS, pages 632-644. Mexico City, Mexico.
- Cristea, D., and Pistol, I. (2009). *The Romanian language from Clarin perspective* (in Romanian). In Diana Trandabăț, Dan Cristea, and Dan Tufiș, editors, Proceedings of the Workshop „Linguistic Resources and Instruments for Romanian Language Processing”, pages 55-64, Iasi, Romania.
- Fillmore, C. J. (1982). *Frame semantics*, în *Linguistics in the Morning Calm*, Hanshin Publishing, Seoul, 111-137.
- Gildea, D. and Jurafsky, D. (2002). Automatic labeling of semantic roles. *Computational Linguistics*, 28(3):245-288.
- Iftene, A., Trandabăț, D., Pistol, I., Moruz, M. A., Husarciuc, M., Cristea, D. (2009). *UAIC Participation at QA@CLEF2008*. In Carol Peters, Thomas Deselaers, Nicola Ferro, Julio Gonzalo, Gareth J.F. Jones, Mikko Kurimo, Thomas Mandl, Anselmo Penas, and Vivien Petras, (eds), *Evaluating Systems for Multilingual and Multimodal Information Access*, 9th CLEF 2008, Revised Selected Papers, vol 5706 of LNCS, pg 385-392. Springer.
- Marquez, L., Carreras, X., Litkowski, K., C. and Stevenson, S. (2008). Semantic role labeling: An introduction to the Special Issue. *Computational Linguistics*, 34(2):145-159.
- Morante, R., Daelemans, W., and Van Asch, V. (2008). *A combined memory-based semantic role labeler of English*. In Proceedings of the Twelfth Conference on Computational Natural Language Learning, pp 208-212, Manchester, UK.
- Nivre, J. (2003) *An efficient algorithm for projective dependency parsing*. In Proc. of the 8th International Workshop on Parsing Technologies (IWPT 03), pp. 149-160.
- Palmer, M., Gildea, D., and Kingsbury, P. (2005). The proposition bank: An annotated corpus of semantic roles. *Computational Linguistics*, 31(1):71-106.
- Postolache, O., Cristea, D. and Orasan, C. (2006). *Transferring Coreference Chains through WordAlignment*. In Proceedings of the 5th Language Resources and Evaluation Conference (LREC2006), Genoa, Italy, May.
- Pradhan, S., Hacioglu, K., Krugler, V., Ward, W., Martin, J. H., and Jurafsky, D. (2005). Support vector learning for semantic argument classification. *Machine Learning Journal*, 60(13):11-39.
- Surdeanu, M., Harabagiu, S., Williams, J., and Aarseth, P. (2003) *Using predicate-argument structures for information extraction*. In Proceedings of the 41th Annual Meeting of the Association for Computational Linguistics, pages 8-15, Tokyo.

- Trandabăț, D. and Husarciuc, M. (2008). *Romanian semantic role resource*. In Proceedings of the Sixth International Language Resources and Evaluation (LREC'08), Marrakech, Morocco, may.
- Trandabăț, D. (2010). *Natural language processing using semantic frames*. PhD Thesis, <http://students.info.uaic.ro/~dtrandabat/thesis.pdf>.
- Trandabăț, D. (2007). *Semantic frames in Romanian natural language processing systems*. In Proceedings of the NAACL-HLT 2007 Doctoral Consortium, pages 29-32, Rochester, New York. Association for Computational Linguistics.
- Tufiș, D., Ion R., Ceașu, Al., Ștefănescu, D. (2005) *Combined Aligners* in Proceeding of the ACL2005 Workshop on „Building and Using Parallel Corpora: Data-driven Machine Translation and Beyond”, Ann Arbor, Michigan, June, 2005
- Tufiș, D., Barbu Mititelu, V., Bozianu, L., and Mihaila, C. (2006). *Romanian wordnet: New developments and applications*. In Proceedings of the 3rd Conference of the Global WordNet Association, pages 337-344, Seogwipo, Jeju, Republic of Korea, Jan.
- Tufiș, D., Ion, R., Ceașu, A., Ștefănescu, D. (2008). RACAI's Linguistic Web Services, în Proceedings of LREC 2008 (Language Resources and Evaluation Conference), May 26 - June 1, Marakkech, Morocco. ELRA - European Language Resources Association
- Witten, I. H. and Frank, E. (2005). *Data Mining: Practical Machine Learning Tools and Techniques (Second Edition)*. Morgan Kaufmann.

SUPRAVEGHEREA PE INTERNET: PĂREREA CONSUMATORILOR DESPRE DIVERSE PRODUSE SAU EVENIMENTE

ADRIAN IFTENE, ALINA-ELENA MIHĂILĂ, GEORGE ALEXANDRU VLAD,
GETA STANCU

Universitatea „Al.I.Cuza”, Facultatea de Informatică, Iași – România

{adiftene, elena.mihaila, george.vlad, geta.stancu}@info.uaic.ro

Rezumat

Lucrarea de față abordează și încearcă să ofere o soluție viabilă la una din problemele zilelor noastre: supravegherea și extragerea de informații relevante de pe Internet. Volumul foarte mare al informațiilor existente pe siteurile web, pe forumuri sau pe pagini personale, fac ca găsirea informațiilor care ne interesează să fie complicată și consumatoare de timp. În plus datorită succesului lor, rețele sociale precum Twitter, MySpace, Facebook, Flickr au un număr de membri cu interese comune în creștere, iar informațiile existente aici au și ele un volum din ce în ce mai mare. Într-o astfel de rețea socială, pornind de la o temă de discuție, utilizatorii își exprimă liber părerile, adaugă link-uri sau poze relevante.

Sistemul pe care l-am construit se folosește de informațiile existente pe Internet și oferă persoanelor care-l folosesc o modalitate ușoară de a afla opinii pozitive sau negative asupra unui subiect de discuție. Informațiile sunt căutate pe paginile web (unde preferăm blog-urile personale, forumurile sau comentariile utilizatorilor) și în rețeaua socială Twitter. Identificarea și clasificarea opiniilor se face prin identificarea unor declanșatori de emoții și prin calcularea unor valențe asociate contextelor din care acestea fac parte.

1. Introducere

În ultimii ani am remarcat o explozie a comunităților de utilizatori în cadrul rețelelor sociale, a forumurilor, sau a site-urilor dedicate unor categorii de persoane care folosesc Internetul pentru a face schimb de mesaje, documente, link-uri, etc. În cadrul acestor comunități, acești participanți își exprimă liber opiniile în legătură cu subiecte de interes comune, critică sau laudă anumite aspecte ale temelor comune de discuție. Interesant este faptul că, în cazul unui produs, de exemplu, pe baza acestor impresii exprimate liber de utilizatori, ne putem face o părere despre calitatea acestuia, putem identifica avantajele folosirii sale și chiar punctele slabe ale acestuia.

De exemplu, site-ul *tweetfeel*¹ vă permite să căutați pe Twitter² ultimele mesaje ale utilizatorilor de aici, în legătură cu un subiect care vă interesează. Aceste postări sunt clasificate în două categorii: prima cu comentarii negative și cea de a doua cu comentarii pozitive pe baza unor cuvinte cheie care sunt identificate în mesajele utilizatorilor. Clasificarea într-una din cele două categorii se face pe baza distanței între cuvintele cheie care caracterizează cele două categorii și cuvintele pe care le căutăm pe

¹ Tweetfeel: <http://www.tweetfeel.com/#iphone>

² Twitter: <http://twitter.com/>

Twitter, fără a se ține cont de contextul în care apar acestea. Deoarece această măsură nu ține cont de semantica textului pentru care se calculează, nu întotdeauna clasificarea comentariilor se face în categoria corectă. Cu toate acestea, contextele în care apar cuvintele cheie folosite de utilizatori în procesul de căutare sunt de ajutor în crearea unei opinii cu privire la ceea ce-l interesează pe acesta.

Lucrarea de față prezintă în prima parte câteva noțiuni legate de ceea ce reprezintă supravegherea utilizatorilor și a informațiilor, și apoi câteva din componentele de bază ale sistemului creat de noi. Aplicația creată are scopul de a afla părerea consumatorilor în legătură cu anumite produse sau evenimente care au avut loc. În cea de a doua parte a lucrării avem prezentate două studii de caz pentru a arăta modul în care funcționează aplicația noastră. În primul, arătăm cum putem afla opinia consumatorilor cu privire la anumite obiceiuri legate de sărbătoarea Paștelui și în cel de al doilea, cum au privit românii, în comparație cu americanii acordarea premiului Nobel președintelui american, Barack Obama.

2. Supravegherea utilizatorilor și a informațiilor de pe Internet

Foucault consideră că în prezent societatea noastră se comportă ca o societate de supraveghere și ca o societate disciplinară. În această societate „*individul este atent fabricat în ea, în conformitate cu o întregă tehnică de forțe și corpuri*” (Foucault, 1977).

Giddens, precizează că la baza supravegherii stă acumularea de informații care pot fi stocate de o agenție sau o colectivitate, precum și supravegherea activităților subordonaților de către superiorii lor (Giddens, 1981). Statul național modern va fi de la început o societate informată, deoarece se ocupă de colectarea și păstrarea informații referitoare la cetățenii săi (nume, adresă, buletin, pașaport, zile de naștere, căsătorii, decese, statistici demografice și fiscale, etc.) în scopul organizării administrării.

Conform (Fuchs, 2009) supravegherea pe Internet este legată de supravegherea informațiile care circulă pe Internet și se face din mai multe motive, din care cele mai importante din ziua de azi sunt legate de securitatea națională (mai ales după atentatul de la 11 Septembrie) și de interese comerciale ale marilor firme. De exemplu, o aplicație care poate extrage informații de pe Internet, legate de acțiunile unor grupuri de utilizatori, care doresc să protesteze cu ocazia unor întâlniri ale oamenilor politici, poate oferi o motivație puternică pentru constituirea unui grup mai mare care să asigure securitatea pe toată durata evenimentului.

Majoritatea tehnicilor ce presupun supravegherea pe Internet implică monitorizarea datelor și a traficului pe Internet. Calculatoarele conectate la Internet comunică între ele cu ajutorul mesajelor, care sunt sparte în bucăți mai mici, pentru a putea fi transmise. Aceste pachete sunt mai apoi transmise prin intermediul unei rețele de calculatoare, din nod în nod, până când își găsesc destinația. Aici, toate pachetele sunt asamblate la loc, formând mesajul inițial. Urmărirea traiectoriei acestor pachete, precum și a conținutului acestora, ține de supravegherea pe Internet.

Conceputul de supraveghere pe internet se află în opoziție cu ideea de securitate a Internetului. Asta deoarece, potrivit Netlingo³ avem că: „*Informațiile care circulă pe Internet, de obicei, au un traseu ocolitor până la calculatorul de destinație, prin mai multe calculatoare intermediare. Traseu real nu este sub controlul celui care dorește să realizeze un astfel de transfer. Prin urmare, la fiecare calculator intermediar există riscul ca cineva să urmărească atent datele care îl traversează și să facă copii ale acestora. Un calculator intermediar poate induce în eroare pe oricine, acesta putând foarte ușor să-și ascundă adevăratele sale intenții. În acest mod acesta poate avea acces la informații confidențiale, cum ar fi parole sau numere de card de credit*”.

3. Prezentarea sistemului

Sistemul construit de noi folosește două modalități de extragere a informațiilor dorite de pe Internet. Prima pornește de la o interfață simplă în care utilizatorul poate introduce o întrebare în limbaj natural. Această întrebare este procesată cu ajutorul componentei de prelucrare a întrebării prezentate în (Iftene et al., 2009) și din ea extragem cuvintele cheie după care face căutarea pe Internet. Această interfață folosind bibliotecile de căutare din Google API⁴ (Google AJAX Search API⁵), extrage de pe Internet legăturile către siteurile cele mai relevante pentru căutarea noastră. În continuare, sistemul nostru apelează la o componentă Lucene Nutch⁶ cu ajutorul căreia salvăm local conținutul siteurilor și apoi le indexează pe calculatorul pe care lucrăm. În final tot cu ajutorul componentei Nutch căutăm în index părerile utilizatorilor, pe care le clasificăm în comentarii pozitive și în comentarii negative folosind un modul specializat în identificarea opiniilor și sentimentelor în texte.

Cea de a doua este similară cu prima, numai că în loc să extragă informațiile de pe Internet folosind Google AJAX Search API, folosește API-ul de la Twitter⁷.

În continuare prezentăm componentele principale ale sistemului, fiecare cu principalele caracteristici.

3.1 Google AJAX Search API

Motoarele de căutare au început să facă parte din ce în ce mai mult din viața noastră de zi cu zi. Volumul foarte mare de informații care există pe Internet, face imposibilă căutare unor informații fără folosirea acestor componente. Folosirea unor astfel de module în cadrul aplicațiilor pe care le construim, duce la mărirea eficienței și a dinamismului sistemului nostru, și la o mai bună interacțiune cu informațiile existente pe Internet.

Pentru a integra în aplicația noastră modulul de căutare Google am folosit bibliotecă JavaScript existente în Google AJAX Search. Acestea ne permit să adăugăm funcționalități noi aplicației noastre, precum ar fi căutarea informațiilor pe Internet.

³ Netlingo: <http://www.netlingo.com/>

⁴ Google API: <http://code.google.com/>

⁵ Google AJAX Search API: <http://code.google.com/apis/ajaxsearch/>

⁶ Lucene Nutch: <http://lucene.apache.org/nutch/>

⁷ Twitter API: <http://apiwiki.twitter.com/>

În urma căutării realizate de utilizator, Google Search API ne întoarce o listă de legături către paginile cele mai relevante. Aceste pagini sunt mai apoi copiate local și indexate cu ajutorul componentelor Nutch.

3.2 Apache Nutch

Nutch/Lucene este o platformă implementată în Java de tip open-source. Componenta Nutch a sistemului nostru este folosită atât pentru copierea locală a siteurilor web cât și pentru indexarea acestora. După procesul de indexare, pe baza unei interfețe web Nutch putem să efectuăm căutări în cadrul indexului creat.

În urma operației de navigare și preluare a informațiilor de pe web sau din cadrul rețelelor locale, conținutul siteurilor este salvat local într-o mulțime de documente denumită *Corpus*. În etapa imediat următoare, acest *Corpus* este indexat și pregătit pentru a putea fi interogată de utilizatori. Ca în cazul celor mai multe motoare de căutare, componenta care se ocupă de partea de interogare este componenta cea mai complicată, dar și cea mai importantă. Interogarea utilizatorului este transformată într-o mulțime de termeni index, care este trimisă mai apoi la motorul de interogare Nutch. Acesta va întoarce documentele cele mai relevante, adică cele în care se găsesc cele mai multe potriviri cu termenii din interogarea utilizatorului.

3.3 Twitter API

Twitter API este bazat în întregime pe protocolul HTTP. Metodele de preluare a datelor din Twitter API necesită o cerere de tip GET, iar metodele care fac introducere, modificare sau distrugere de date necesită cereri de tip POST. Sunt permise de asemenea cereri de tip DELETE, atunci când dorim să distrugem anumite informații.

Twitter API permite modificarea formatului extensiei cererii pentru a obține rezultatele într-un alt format dorit de utilizator. Din formatele suportate în prezent de acest API (XML, JSON și RSS sau Atom) aplicația construită de noi a folosit formatul XML.

Din funcțiile existente în acest API am folosit în primul rând funcțiile de căutare, care ne permit să extragem informațiile dorite din postările utilizatorilor.

3.4 Identificarea opiniilor pozitive și negative ale utilizatorilor

Identificarea înțelesului cuvintelor din texte a devenit în ultimii ani o direcție importantă de cercetare în domeniul lingvisticii computaționale. Abordările existente folosesc dezambiguizare, folosesc resurse lingvistice precum SensiWordNet sau încearcă identificarea declanșatorilor de emoții (vezi lucrările (Hatzivassiloglou și Mckeown, 1997), (Mihalcea et al., 2007) și (Tufiş, 2009)).

Pentru identificarea opiniilor am folosit o metodă similară celei folosite în (Iftene și Rotaru, 2010), care construiește incremental o bază de date lexicală (care conține cuvinte ce declanșează emoții) similar cu (Balahur și Montoyo, 2008). Elementele din această bază de date ne ajută să descoperim opiniile și emoțiile în textele extrase de API-urile Google și Twitter, iar în urma calculării unor valențe globale ale textelor ne putem da seama dacă acestea reprezintă latura pozitivă, negativă sau neutră a unui sentiment.

Elementul de bază care identifică emoțiile este denumit *declanșator de emoții* și reprezintă un cuvânt sau un concept care poate oferi o interpretare emoțională a conținutului textului. Iată câteva exemple de declanșatoare de emoții: cuvinte precum „mândrie”, „libertate”, „stimă”, „familie”.

Baza lexicală de termeni a fost construită pornind de la 30 de termeni prezenți în „piramida lui Maslow” (Maslow, 1943), care au fost traduși în limba română. Adicional am folosit WordNet-ul românesc (Tufiș et al., 2004) de unde am preluat sinonimele, antonimele, și hiponimele.

După construirea bazei lexicale de termeni, următorul pas a fost de a atribui valențe și emoții termenilor din această bază de date. Pentru aceasta am ținut cont de următoarele reguli:

- Termenilor principali de declanșare a emoțiilor și sinonimelor acestora li se acordă o valoare pozitivă.
- Termenilor hiponimi cu cei de mai sus și termenilor derivați din termeni principali care sunt verbe li se acordă de asemenea o valoare pozitivă.
- Termenilor antonimi cu oricare din cei de mai sus li se acordă o valoare negativă.
- Valențele oricăror termeni se modifică în funcție de modificatorii (termeni care neagă, accentuează sau diminuează o valență) care îi însoțesc.

3.4.1 Modificatori de valență

Pentru a determina valențele finale ale declanșatoarelor de emoții se definește o mulțime de *modificatori de valență* (Balahur și Montoyo, 2008). Un **modificator de valență** reprezintă un termen care modifică valența asociată unui alt termen pe care îl însoțește. Modificatorii pe care i-am identificat sunt de mai multe feluri:

- **Negații** – care modifică valența radical de la polul pozitiv la cel negativ sau invers. Aici modificatorii sunt reprezentați de cuvinte care introduc negația: „nu”, „niciodată”, etc.
- **Accentuatori** – care accentuează aspectul negativ sau pozitiv al unui declanșator. Din această mulțime fac parte adjective precum: „mare”, „mult”, „bine”, „profund”, „excepțional” sau adverbe care accentuează înțelesul întregului context din care fac parte „cu siguranță”, „sigur”, „cert”, „în definitiv”.
- **Diminuatori** – care diminuează aspectul negativ sau pozitiv al unui declanșator ducându-l spre o valență neutră. Diminuatorii sunt reprezentați de adjective precum „mic”, „puțin”, „rău”, „degrabă”, de verbe modale: „a putea”, „a fi posibil”, „a trebui”, „a vrea”, de adverbe ca „posibil”, „probabil”. Verbele modale și adverbele introduc noțiunea de incertitudine și posibilitate și diminuează valența emoției întregului context din care fac parte. De asemenea acești diminuatori ne ajută să facem distincția între evenimente care au avut loc, ar fi putut să aibă loc, au loc în prezent sau vor avea loc în viitor.

De exemplu, cum ar suna un fragment în care există modificatori de valență diferiți: „*este minunat*” (forma inițială), „*nu este minunat*” (forma negată), „*ar putea fi minunat*” (forma diminuată), „*e absolut minunat*” (formă accentuată).

Există însă și termeni care, dacă sunt combinați/alăturați, produc sentimentul de ironie. Acesta este și cazul asocierii următoare: „*Exceptionalul organizator a eșuat în rezolvarea problemei.*”, unde termenul „*exceptional*” (are în mod uzual valență pozitivă) alături de termenul „*a eșuat*” (are valență negativă de obicei) ne duce cu gândul la faptul că „*organizatorul care era recunoscut pentru priceperea sa a dat-o în bară*”.

În plus, pe parcursul testării și verificării aplicației, am identificat termeni noi specifici discuțiilor purtate pe forumuri, blog-uri sau rețele sociale și cu ajutorul lor am completat resursele obținute în pașii anteriori. Iată câteva exemple de termeni specifici: „*brio*”, „*super*”, „*fine*”, „*bună*” sau succesiuni de caractere speciale care reprezintă iconițe emoționale: „*:)*” (față zâmbitoare ☺), „*:(*” (față tristă ☹), etc.

O observație importantă este următoarea: pentru ca un modificador de valență să își îndeplinească scopul (să modifice valența termenilor), trebuie ca în text să fi fost exprimată o atitudine (al cărei înțeles să poată fi modificat). De exemplu, în propoziția „*Ion este acasă.*” care nu face decât să prezinte un fapt și nu o atitudine, introducerea unei negații „*Ion nu este acasă*” nu schimbă valența nici unui termen.

4. Studii de caz

În acest capitol am realizat două studii de caz, pentru a arăta modul de funcționare al sistemului nostru. Cu ajutorul acestui sistem utilizatorii vor avea toate aceste informații centralizate și în plus ele vor fi clasificate în comentarii pozitive și negative. Nu în ultimul rând aplicația, prin faptul că accesează simultan diferite zone de pe Internet și că oferă informațiile centralizat, îi va ajuta pe cei care o folosesc să reducă timpul necesar obținerii unor astfel de informații.

4.1 Sărbătoarea Paștelui

Primul studiu de caz este legat de sărbătoarea Paștelui. Scopul aplicației este de a veni în ajutorul celor care doresc să realizeze produse specifice acestei sărbători, ei neavând la dispoziție rețete disponibile și nici timp pentru a le căuta. Pentru acest studiu de caz am luat în considerare produsele tradiționale de Paște (cozonac, pască și ouă roșii) și am identificat opiniile utilizatorilor de pe Internet cu privire la anumite rețete sau modalități de preparare. Scopul final a fost de a identifica dacă putem găsi cu ajutorul aplicației o rețetă cât mai apropiată de gustul utilizatorului și care sunt avantajele și dezavantajele acesteia.

Pentru API-ul de la Google, pentru a fi siguri ca rezultatele întoarse vor fi relevante, am construit diverse interogări alegând diverse combinații ce folosesc cuvinte care reprezintă fie produse (*cozonac, pască, ouă roșii*), fie ingrediente ale acestora (*ciocolată, brânză, smântână, nucă*) (cu sau fără diacritice) și în plus am cerut ca aceste informații să se găsească pe pagini ce respectă formate de *forum* sau *blog*.

Pentru API-ul de Twitter am construit interogări mai simple (cu mai puține cuvinte) pentru a avea șanse mai mari de reușită folosind din nou cuvinte simple precum cele de mai sus. În final am realizat peste 20 de căutări pe Google și 7 pe Twitter.

Pentru fiecare căutare realizată pe Google, Nutch-ul a preluat maxim 10 link-uri din rezultatele oferite pentru care a salvat local conținutul site-urilor. Acest conținut l-am procesat folosind componenta de identificare a emoțiilor și am extras sentimentele pozitive și negative.

În continuare vom vedea câteva exemple de propoziții extrase cu ajutorul API-ului Google de pe forumuri în care am îngroșat cuvintele pentru care am identificat valențe pozitive sau negative. Astfel, în urma căutărilor realizate pe tematici cu rețete de cozonaci, în comentariile utilizatorilor am găsit cuvinte declanșatoare de emoții cu tentă pozitivă „se **pare** ca cozonacul mamei a trecut cu **brio** proba :)”, „Aluatul de cozonac e **foarte delicat și sensibil.**”, „E o rețetă **super** simplă și **super** gustoasă...”, și cu tentă negativă „oricât aș vrea **nu** reușesc sa fac acest cozonac”, etc.

Cu toate că în Twitter căutările au fost mult mai simple, rezultatele au fost considerabil mai puține decât cele întoarse de motorul Google, deoarece am observat că numărul comentariilor în limba română sunt foarte puține și de cele mai multe ori conțin cuvinte și chiar fragmente în limba engleză.

Iată câteva mesaje obținute de pe Twitter în urmă căutărilor realizate pentru *cozonac*, *pască* și *ouă roșii*: „deza@zozo_ro: **ciudată** combinație...dar suna gustos...acum mănânc din primul meu cozonac...si **ciudat**, a ieșit **bun**”, „richieTM: Deci cea mai **buna** plăcinta cu brânza e Pasca!!”, „lau_anca @RaduCeuca **super! dar** ai uitat de clasicele oua roșii :)”, etc.

În urma acestui studiu am observat că opiniile utilizatorilor de forum-uri și blog-uri sunt mai consistente și mai relevante decât cele ale utilizatorilor de pe Twitter, care sunt mai scurte și care se folosesc de multe ori de imagini încărcate de aceștia.

4.2 Acordarea premiului Nobel lui Obama

Pe data de 10 octombrie, știrea potrivit căreia Barack Obama a câștigat Premiul Nobel pentru Pace a ținut prima pagină a ziarelor, atât în România cât și în străinătate. Studiul de caz de față își propune să analizeze opiniile oamenilor vorbitori de limba română (și pentru aceasta am căutat pe site-urile, blog-urile și forum-urile scrise în limba română folosind interogări care folosesc cuvinte românești) și a celor vorbitori de limba engleză (cu ajutorul interogărilor scrise în limba engleză) cu privire la aceasta decizie.

Rezultatele obținute au arătat clar o diferență între comentariile celor din America (unde comentariile au fost în mare parte pozitive) și celor din România (unde comentariile au fost în mare parte negative). Trebuie să precizăm că în America părerile au fost exprimate de persoane care au avut un rol important în alegerea lui Obama ca președinte, la mai puțin de un an de când aceștia l-au votat datorită discursurilor pe care acesta le-a ținut în timpul campaniei sale. În același timp în România, părerile au fost exprimate de persoane care nu au fost implicate în nici un fel în alegerile din America.

Aplicația care folosește API-ul Google a extras informații de pe 20 de sururi (zece din România și zece din străinătate), însumând în total 327 de comentarii. Pentru a

identifica valențele fragmentelor extrase pentru vorbitorii de limba engleză am tradus folosind Google Translate⁸ resursele obținute pentru limba română și le-am extins folosind WordNet-ul⁹ englezesc.

Cele 20 de situri au fost obținute în urma folosirii unor interogări ce foloseau cuvinte cheie corespunzătoare unor întrebări ca mai jos:

1. Ce părere aveți despre înmânarea premiului Nobel lui Barack Obama? (En: *What do you think about giving the Nobel Prize to Barack Obama?*)
2. Care sunt părerile pro și contra la acordarea premiului Nobel lui Obama? (En: *What are the pros and cons to award the Nobel Prize Obama?*)
3. A meritat Obama premiul Nobel? (En: *Did Obama deserved the Nobel Prize?*)

În ceea ce privește siturile din România, o mare parte din cei care și-au exprimat părerea au scos în evidență faptul că datorită acestei alegeri neinspirate premiile vor avea de suferit de acum încolo. Odată cu luarea acestei decizii, mulți cred că premiile și-au pierdut din importanță și că de acum câștigarea lor nu va mai fi o realizare atât de mare, întrucât un precedent a fost deja creat.

Iată câteva exemple de comentarii adăugate de utilizatori, în care am detectat o mare încărcătură emoțională corespunzătoare dezaprobării și indignării (toate având aceeași sursă: http://economie.hotnews.ro/stiri-media_publicitate-6262253-presa-americana-radiosul-obama-accepta-premiul-pentru-pace-sau-cum-castigi-nobelul-12-zile.htm):

- „...e **culmea tupeului** sa dai premiul Nobel cuiva care era presedinte de doar 12 zile la data inchiderii perioadei de candidatura !!...”
- „acest Nobel al lui Obama va ramane in istorie ca fiind **probabil** cel mai **nemeritat** din cate s-au acordat!”
- „**Cred** ca prin aceasta **trista** alegere, nemeritata, Nobelul pt Pace e definitiv **compromis**.” , „a devenit evident ca **nu** are vreo legatura cu meritele sau rezultate reale obtinute. ...”
- „**Cred** ca e o **jignire** la adresa celor care chiar fac ceva pentru pace in lume...” , „...pentru ce?!?! Pentru cele doua **razboaie** in care SUA sunt implicate?”
- „Un act de o **nesimtire** și un **dispret incredibil** fata de sute de lideri politici”.

În siteurile în engleză (în principal pentru cele din America) am observat o altă atitudine. Deși marea majoritate nu sunt de acord cu primirea acestei distincții, criticile nu sunt îndreptate direct către președintele American, ci mai mult spre comisia care a acordat această distincție.

În legătură cu declanșatoarele de emoții, siturile în limba română conțin multe cuvinte cu valențe negative și modificatori care întăresc aceste valențe (precum *incredibil*, *total*, *foarte*). În limba engleză se observă o tendință de a folosi cuvinte ce neutralizează valențele (*poate că*, *posibil*, *probabil*).

⁸ Google Translate: <http://translate.google.com/>

⁹ WordNet: <http://wordnetweb.princeton.edu/perl/webwn>

În multe reacții am observat o nedumerire generală („nu înțeleg ce treaba are”, „nu văd de ce”, „nu știu de ce a fost ales”, „nu cred că trebuia...”, „nu consider că aceasta era cea mai bună alegere”) cu privire la luarea acestei decizii.

Puținii oameni care nu comentează negativ această știre, cu toate că nu sunt neapărat de acord cu luarea acestei decizii, își exprimă o admirație față de persoana în cauză.

5. Concluzii

Această lucrare prezintă principalele componente ale sistemului pe care-l oferim utilizatorilor care doresc să afle opinii pozitive sau negative asupra unor produse sau evenimente. Sistemul este o alternativă la căutarea clasică pe Internet, aducând nou posibilitatea de combinare a rezultatelor obținute în urma căutărilor uzuale cu rezultate extrase de pe blog-urile utilizatorilor, forumuri sau chiar rețele sociale. Componentele principale ale sistemului se bazează pe componentele care ne ajută să căutăm și să extragem informații relevante de pe Internet: API-urile de căutare de la Google și de la Twitter.

Sistemul folosit de noi a folosit în calculul valențelor asociate propozițiilor distanța între cuvintele cheie și declanșatorii de emoții, ținând cont de modificatorii de valență. Din evaluările făcute, pe 6 teme de căutare efectuate, cu ajutorul a 3 persoane care au evaluat peste 100 de paragrafe extrase putem spune că suntem la început de drum, calitatea rezultatelor fiind încă modestă (în jur de 44%). Principalele probleme sunt datorate în principal faptului că declanșatorii de emoții vizați nu se refereau la cuvintele cheie pentru care realizăm căutarea. O altă problemă importantă este legată de faptul că prin modul de calcul al valențelor asociate, preferăm propozițiile scurte, care nu sunt întotdeauna relevante pentru ceea ce dorim să obținem.

În continuare, dorim să implementăm în modul de calcul al valențelor asociate elemente care să elimine principalele probleme prezentate mai sus, bazate pe identificarea rolurilor semantice. De asemenea dorim să realizăm o evaluare mai relevantă, cu circa 20 de teme de căutare și cu mai mult de 10 evaluatori umani.

Mulțumiri. Cercetarea prezentată în această lucrare a fost finanțată de către Programul Operațional Sectorial Dezvoltarea Resurselor Umane prin proiectul „Dezvoltarea capacității de inovare și creșterea impactului cercetării prin programe post-doctorale POSDRU/89/1.5/S/49944.

Referințe bibliografice

- Balahur, A., Montoyo, A. (2008). *Applying a culture dependent emotion triggers database for text valence and emotion classification*. In journal Procesamiento del Lenguaje Natural, ISSN 1135-5948, N°. 40. Pp. 107-114.
- Foucault, M. (1977). *Discipline and punish*. New York: Vintage.
- Fuchs, C. (2009). *Social Networking Sites and the Surveillance Society*. A Critical Case Study of the Usage of studiVZ, Facebook, and MySpace by Students in Salzburg in the Context of Electronic Surveillance Salzburg/Vienna, Austria. 2009.

- Forschungsgruppe „Unified Theory of Information” - Verein zur Förderung der Integration der Informationswissenschaften. ISBN 978-3-200-01428-2.
- Giddens, A. (1981). *A contemporary critique of Historical Materialism*. Vol. 1: Power, property and the state. London: Macmillan.
- Hatzivassiloglou, V., Mckeown, K. R. (1997). *Predicting the semantic orientation of adjectives*. In Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics, Morristown, NJ, USA. pp. 174-181.
- Iftene, A., Rotaru, A. (2010). *User Profile Modeling in eLearning using Sentiment Extraction from Text*. In Research in Computing Science, „Special issue: Natural Language Processing and its Applications”, Vol.46, Pp.267-278, Instituto Politecnico Nacional, Centro de Investigacion en Computacion, Mexico 2010. ISSN: 1870-4069. Poster at 11th International Conference on Intelligent Text Processing and Computational Linguistics (CICling 2010). 21-27 March. Iasi, Romania.
- Iftene, A., Trandabăț, D., Pistol, I., Moruz, A.M., Husarciuc, M., Sterpu, M. and Turliuc, C. (2009). *Question Answering on English and Romanian Languages*. In Proceedings of the CLEF 2009 Workshop. 30 September - 2 October. Corfu, Greece.
- Maslow, A. H. (1943). *A Theory of Human Motivation*. Psychological Review 50 (4). 370-96.
- Mihalcea, R., Banea, C. și Wiebe, J. (2007). *Learning Multilingual Subjective Language via Cross-Lingual Projections*. In Proceedings of the Association for Computational Linguistics (ACL 2007), Prague, June.
- Tufiș, D. (2009). *Playing with Word Meanings*. In Lotfi A. Zadeh, Dan Tufiș, Florin Gh. Filip and Ioan Dzițăc (eds.) *From Natural Language to Soft Computing: New Paradigms in Artificial Intelligence*, Editing House of Romanian Academy. ISBN 978-973-27-1678-6, pp. 211-223.
- Tufiș, D., Barbu, E., Barbu Mititelu, V., Ion, R., Bozianu, L. (2004). *The Romanian Wordnet*. Romanian Journal of Information Science and Technology, Volume 7, Numbers 1-2, pp. 107-124.

FOLOSIREA VERBELOR PENTRU DETERMINAREA INFERENȚELOR TEXTUALE

MIHAI ALEX MORUZ^{1,2}

¹Universitatea „Al.I.Cuza”, Facultatea de Informatică, Iași – România

²Academia Română, Filiala Iași, Institutul de Informatică Teoretică – România

mmoruz@info.uaic.ro

Rezumat

Această lucrare descrie o metodă de utilizare a verbelor și a structurii de argumente a acestora în vederea rezolvării de inferențe textuale pentru texte în limbile engleză și română. Analiza efectuată asupra unui set de perechi de inferențe din setul de test folosit la RTE-5 dovedește că semantica predicțională, așa cum a fost definită de Charles Fillmore și Beth Levin, este utilă pentru rezolvarea inferențelor textuale, întrucât 38% dintre perechile analizate sunt rezolvate direct folosind acest tip de semantică, iar 11% dintre perechi sunt rezolvate recurgând la elemente analiză verbală.

1. Introducere

Noțiunea de inferențe textuale (Textual Entailment în limba engleză) a fost descrisă de (Dagan, Glickman, 2004) ca o metodă de a generaliza noțiunea de variabilitate a limbajului natural în diverse probleme din domeniul limbajului natural, cum ar fi sistemele întrebare – răspuns, sumarizarea automată, etc. Conform (Dagan, Glickman, 2004), inferențele textuale sunt definite ca relație dintre un text coerent, **T**, și o expresie în limbaj natural, care este văzută ca o ipoteză, **H**. Spunem că **T** inferă **H**, (**H** este o consecință a lui **T**), notat $T \Rightarrow H$, dacă înțelesul lui **H**, interpretat în contextul lui **T**, poate fi dedus din înțelesul lui **T**. Relația de inferență definită mai sus este direcțională, deoarece deși un text poate fi dedus din un altul, reciproca nu este întotdeauna valabilă.

Definiția de mai sus, deși completă și corectă, este prea abstractă pentru a putea fi folosită direct în sisteme de procesare a limbajului natural. Recognising Textual Entailment Challenge (RTE) este principala campanie de evaluări pentru sistemele de rezolvare a inferențelor textuale, și a fost inițiată în 2005 de către PASCAL (Pattern Analysis, Statistical Modelling and Computational Learning – <http://www.pascal-network.org/>) (Dagan et al., 2005), iar ultimele două ediții, cele din 2008 și 2009, au fost organizate în cadrul Text Analysis Conference (TAC – <http://www.nist.gov/tac/>). Din acest motiv au apărut o serie de adaptări ale definiției inițiale pentru inferențele textuale, care să poate fi utilizate în abordări practice.

Analizând diverse seturi de perechi de inferență, am ajuns la intuiția că majoritatea cazurilor pot fi rezolvate prin examinarea a două tipuri de informație:

- Relația dintre verbele din ipoteză și cele din text. Verbul trebuie considerat împreună cu argumentele și adjuncții săi; așadar, compararea dintre două verbe implică o comparare între două structuri complexe care sunt, în esență, propoziții atomice (clauze). Dacă verbele din **H**, împreună cu argumentele și adjuncții lor, sunt

susținute în **T**, rezultatul este ENTAILMENT; acest lucru se face prin identificarea acelor elemente din **T** care suportă informația din **H**.

- Fiecare argument sau adjunct este o entitate, care are atașat un set de proprietăți. Este posibil ca, în ciuda potrivirii la nivel de verbe, să apară diferențe în ceea ce privește proprietățile argumentelor sau adjunctilor similari. Pentru a putea rezolva aceste cazuri, considerăm fiecare adjunct sau argument o entitate care are un set de atribute atașat, iar compararea se face atât pe nucleul entității cât și pe setul de atribute.

Această lucrare are ca scop principal determinarea relevanței verbelor pentru problema inferențelor textuale, așa cum au fost ele definite în cadrul acestei secțiuni. Mulțimile de proprietăți atașate argumentelor sau adjunctilor sunt discutate pe scurt în secțiunea 4.3, cu referire la dezvoltări ulterioare.

Ideea folosirii verbelor pentru a rezolva inferențele textuale nu este nouă, și a fost utilizată de (Burchardt, Frank, 2006) și (Hickl, Bensley, 2007), (Hickl, 2008), printre alții. (Burchardt, Frank, 2006) aproximează inferențele textuale cu suprapunerea la nivel structural și semantic între text și ipoteză, și combină parsarea LFG cu semantica cadrelor, pentru a proiecta o reprezentare lexical semantică cu roluri semantice, și apoi calculează suprapunerea dintre **H** și **T**. (Hickl, Bensley, 2007), (Hickl, 2008) extrag o mulțime de propoziții acceptate universal, numite contexte discursive (*discourse commitments*); odată ce acestea au fost extrase, problema rezolvării inferențelor textuale se reduce la determinarea acordului contextelor discursive din **H** cu cele din **T**.

Diferența esențială dintre abordările precedente și soluția descrisă în cadrul acestei lucrări este aceea că soluția noastră folosește noțiunea de alternanțe verbale și rezultatele lui Beth Levin asupra claselor verbale (cunoscute drept *clase Levin*, și descrise în (Levin, 1993)). Mai precis, este determinată clasa Levin pentru fiecare verb din **H**, iar apoi se încearcă găsirea acelor verbe din **T** care au aceeași clasă Levin și care susțin ipoteza; dacă aceasta eșuează, se încearcă determinarea relațiilor dintre verbe pe baza descrierilor semantice atașate fiecărei clase. Această procedură a fost testată prin analiza detaliată a unui set de 200 de perechi de inferențe din datele de test din RTE-5¹, și am determinat că peste 38% dintre ele pot fi rezolvate corect folosind clase Levin. Pentru identificarea claselor Levin am folosit VerbNet², descris în detaliu în secțiunea 2.

O altă contribuție nouă a acestei abordări este gruparea proprietăților entităților în seturi de proprietăți; pentru două entități similare în **T** și **H**, inferența poate fi definită ca rezultatul comparării acestor seturi de atribute: ENTAILMENT înseamnă ca toate proprietățile din **H** sunt subsumate de proprietăți din **T**. Din cele 200 de perechi analizate, peste 29% pot fi rezolvate pe baza acestei metode.

2. Clasele Levin și VerbNet

Cea mai mare și mai frecvent utilizată clasificare a verbelor din limba engleză este cea descrisă în (Levin, 1993), mai larg cunoscută drept clase Levin. VerbNet (VN) (Kipper

¹ TAC 2009 Recognizing Textual Entailment Track, <http://www.nist.gov/tac/2009/RTE/>

² <http://verbs.colorado.edu/~mpalmer/projects/verbnet.html>

et al., 2000; Kipper-Schuler, 2005) este un lexicon de verbe din limba engleză disponibil online care oferă descrieri sintactice și semantice detaliate pentru clasele Levin, care sunt organizate într-o taxonomie mai rafinată. Conform (Kipper-Schuler et al., 2006) VN este un lexicon de verbe ierarhic, independent de domeniu, cu acoperire largă, care este aliniat cu resurse ca FrameNet (Baker et al., 1998) și WordNet (Fellbaum 1998).

O clasă VerbNet este complet descrisă de o mulțime de verbe membre, rolurile tematice pentru structurile predicat-argument ale acestor verbe, restricțiile de selecție pentru argumente și cadre care conțin o descriere sintactică și predicate semantice care au și o funcție temporală, similar cu descompunerea de evenimente propusă de (Moens, Steedman, 1988). VerbNet extinde clasele Levin inițiale prin rafinarea acestora în subclase cu un mai mare grad de coerență sintactică și semantică. Această ierarhie a fost extinsă, în urma studiului efectuat de (Korhonen, Briscoe, 2004), cu 57 de noi clase (de la aproximativ 200) și 106 de noi alternanțe de diateză (de la aproximativ 200).

Întrucât această resursă este aliniată cu WN, ea poate fi ușor transferată pentru limba română, prin intermediul synseturilor, ceea ce ar permite aplicarea metodei descrisă în secțiunile 3 și 4 pentru rezolvarea inferențelor textuale în limba română.

2.1 Cadrele sintactice în VN

Fiecare clasă VN conține un set de descrieri sintactice (cadre) care reprezintă posibile realizări de suprafață pentru structurile de argumente (tranzitive, intransitive, cu grupuri prepoziționale etc.) și un număr de alternanțe de diateză, date de Levin ca parte a fiecărei clase (Kipper-Schuler, 2006). Fiecare cadru sintactic este descris folosind roluri tematice (*Agent, Patient* etc.), verbul și toate celelalte elemente lexicale necesare pentru a construcție sau alternanță. Rolurile tematice pot fi restricționate semantic (*animate, organization, human* etc.), iar cadrele sintactice pot fi și ele restricționate cu privire la prepozițiile permise. Restricții adiționale pot fi impuse asupra rolurilor tematice, pentru a transmite natura sintactică a constituentului care este cel mai probabil asociat cu rolul tematic (*sentential, nominal* etc).

2.2 Predicate semantice

În afara descrierilor sintactice atașate, fiecare clasă Levin are și o descriere semantică, în forma unei conjuncții de predicate semantice booleene precum *cause, motion, sau contact*. Fiecare dintre aceste predicate este asociat cu o variabilă eveniment *E*, care este folosită pentru specificarea momentului la care predicatul este adevărat. Relațiile dintre verbe sau dintre clase de verbe (antonimie, consecință, etc) pot fi deduse pe baza acestor predicate semantice. În VN, aspectul unui verb este descris prin intermediul variabilei eveniment din predicate (Kipper-Schuler et al., 2006).

2.3 Statistici legate de VerbNet

În acest moment, VN versiunea 3.1, cea mai recentă versiune de VN, conține descrieri pentru peste 5800 de sensuri verbale, distribuite în 270 de clase de nivelul întâi și 200 de subclase. Descrierile acestor verbe folosesc 33 de roluri tematice, 36 de restricții de selecție, 296 de cadre principale (care pot fi extinse prin specificarea tipurilor de roluri tematice sau tipuri de complemente) și 145 de predicate semantice. Lexiconul se

bazează și pe o ierarhie de prepoziții cu 66 de elemente. Acoperirea PropBank (Palmer et al., 2005) de către VN a crescut în această versiune la peste 90%. În plus, în VN 3.1 taxonomia definită de (Levin, 1993) a fost extinsă în profunzime, prin rafinarea unor clase existente și în lărgime, prin adăugarea claselor propuse de (Korhonen, Briscoe, 2004).

3. Algoritmul preliminar pentru compararea verbelor

Scopul studiului de fezabilitate, care va fi descris în detaliu în secțiunea 4, este acela de a testa dacă rezolvarea inferențelor textuale folosind relații între verbe, definite pe baza claselor Levin, este posibilă; rezultatele au arătat că acest fapt nu este doar posibil, ci că o astfel de abordare rezolvă direct peste 38% dintre perechile luate în considerare, și ajută, indirect, la rezolvarea a peste 11% altor perechi (o descriere detaliată este dată în secțiunea 4). Din acest motiv propunem un algoritm preliminar pentru rezolvarea inferențelor textuale pe baza relațiilor dintre verbe.

Algoritmul primește la intrare arbori sintactici pentru **T** și **H**, și întoarce rezultatul inferenței pentru perechea dată. Descrierea preliminară a algoritmului pentru rezolvarea de inferențe textuale pe baza relației dintre verbe este dată mai jos.

1. Se extrag clasele Levin pentru toate verbele din **T** și **H** și se atașează descrierea semantică corespunzătoare.
2. Se determină dacă fiecare verb din **H** are aceeași clasă Levin cu măcar un verb din **T**
 - 2.a. Dacă argumentele și adjunctii se potrivesc peste **T** și **H**, iar verbele nu sunt antonime, ENTAILMENT
 - 2.b. Altfel, dacă argumentele și adjunctii se potrivesc dar verbele sunt antonime, sau dacă argumentele nu se potrivesc (sunt entități conflictuale sau aceeași entitate dar cu atribute conflictuale), CONTRADICTION
 - 2.c. Altfel, UNKNOWN
3. Dacă nu există similaritate de clase Levin, se extrag relațiile dintre verbe pe baza descrierii semantice.
 - 3.a. Dacă verbele din **H** sunt sinonime (au descrieri semantice similare conform cadrelor VN) sau consecințe (starea sau acțiunea descrisă în **H** este un rezultat al evenimentului din **T**) ale verbelor din **T**, iar argumentele se potrivesc, sau dacă descrierea semantică a verbului corespunde cu o pereche element-atribut din **T**, ENTAILMENT
 - 3.b. Altfel, dacă verbele din **H** sunt antonime cu verbe din **T** sau mai generale decât acestea, iar argumentele se potrivesc, CONTRADICTION
 - 3.c. Altfel, UNKNOWN

Algoritmul descris mai sus folosește clasele Levin, așa cum sunt descrise în VN, pentru a rezolva inferențe textuale. Primul pas este acela de a atașa clasele Levin

corespunzătoare pentru fiecare verb din **T** și **H**. Apoi încearcă extragerea de propoziții atomice echivalente pe baza claselor Levin. Dacă verbele astfel găsite se potrivesc la nivelul argumentelor și adjuncțiilor, soluția este ENTAILMENT; dacă argumentele nu se potrivesc sau verbele sunt în aceeași clasă dar sunt antonime, soluția este CONTRADICTION; în orice alt caz, rezultatul este UNKNOWN. Dacă nu apar potriviri la nivel de clase Levin, algoritmul încearcă extragerea de propoziții atomice echivalente pe baza descrierii semantice a verbelor, iar apoi face pași similari ca mai sus. Dacă există verbe în **H** care nu au echivalent în **T** în urma pașilor descriși mai sus, rezultatul este UNKNOWN.

4. Analiza corpusului

Așa cum am afirmat și în introducere, una dintre metodele cheie pentru determinarea inferențelor textuale, conform intuiției noastre, este analiza verbelor. Scopul acestei secțiuni este acela de a prezenta studiul de fezabilitate efectuat pentru a determina utilitatea claselor Levin pentru problema inferențelor textuale.

Testul de fezabilitate a fost realizat asupra unui corpus semnificativ, de 200 de perechi de inferență, extrase din setul de test RTE-5 (acesta este cel mai recent corpus de inferențe textuale pentru limba engleză). Setul de perechi ales este tipic pentru datele RTE în ceea ce privește distribuția perechilor: dintre cele 200 de perechi analizate, 100 (50%) sunt ENTAILMENT, 31 (15.5%) sunt CONTRADICTION și 69 (35.5%) sunt UNKNOWN.

În urma analizei, fiecare pereche a fost clasificată în una dintre categoriile de mai jos:

- **VN (VerbNet)**, ceea ce înseamnă că aplicarea algoritmului definit în secțiunea 3 conduce la determinarea soluției corecte pentru perechea dată, în urma faptului că verbele din **T** și **H** au aceeași clasă Levin, au descrieri semantice similare, sunt consecințe pentru verbe din **T** sau pe baza compatibilității structurii de argumente;
- **NA (NeAtașate)**, ceea ce înseamnă că algoritmul din secțiunea 3 nu poate găsi un verb în **T** care să aibă aceeași structură de argumente ca un verb din **H**, sau nu există verbe în **T** care să fie conectate semantic de verbe din **H**;
- **L (Lipsă)**, ceea ce înseamnă că unul dintre conceptele cheie din **H** nu apare în **T**;
- **AS (Alte Surse)**, ceea ce înseamnă că rezultatul corect poate fi dedus prin alte mijloace decât cele date mai sus (de exemplu rezoluția anaferei sau cunoaștere ontologică). Este important de menționat că, în urma analizei noastre, cele mai multe dintre perechile care cad în această categorie pot fi reduse la cazuri de comparare de seturi de trăsături pentru argumente și adjuncți.

Tabelul 1: Distribuția soluțiilor pentru studiul de fezabilitate

Categoria	Număr	Procent
VN	76	38%
AS	59	29.5%
L	43	21.5%
NA	22	11%

După cum se poate vedea din distribuția soluțiilor pentru perechile analizate, cel mai mare câștig vine din folosirea VerbNet, urmat de alte surse de inferență (aproape toate aceste cazuri sunt instanțe de comparare de seturi de proprietăți). Distribuția susține intuiția noastră, întrucât toate cazurile de **NA** pot fi reduse la predicate care nu sunt echivalente, iar toate cazurile de **L** pot fi reduse la seturi de proprietăți care nu se potrivesc.

4.1 Perechi rezolvate folosind VN

Cea mai simplă utilizare a VN este alinierea verbelor din cadrul aceleiași clase Levin, așa cum se vede în exemplul 1.

Ex. 1: Potrivire exactă peste clase VN

T: *MADAGASCAR'S constitutional court declared Andry Rajoelina as the new president of the vast Indian Ocean island today... (Curtea constituțională din Madagascar l-a declarat pe Andry Rajoelina președinte al marii insule din oceanul Indian)*

H: *Andry Rajoelina was proclaimed president of Madagascar. (Andry Rajoelina a fost proclamat președinte al Madagascarului.)*

Verbul *proclaim* din ipoteză face parte din clasa *say-37.7-1*, care conține și verbul *declare* din text. Aceste două verbe se potrivesc din punct de vedere al argumentelor (*Andry Rajoelina* este *Theme*) și al adjuncțiilor (*president of Madagascar*). Soluția pentru această pereche este ENTAILMENT, conform pasului **2.a.** din.

Folosirea descrierii sintactice a unui cadru și a descompunerii semantice atașate este descrisă în exemplul 2.

Ex. 2: Descriere sintactică și descompunere semantică

T: *A court in Venezuela has jailed nine former police officers for their role in the deaths of 19 people during demonstrations in 2002. ... (Un tribunal din Venezuela a încarcerat nouă foști polițiști pentru rolul lor în morțile a 19 persoane în timpul demonstrațiilor din 2002)*

H: *Nine police officers have had a role in the death of 19 people. (Nouă polițiști au avut un rol în moartea a 19 persoane.)*

Verbul *have* face parte din clasa *own-100*, care este descrisă sintactic *Theme1 V Theme2* și semantic *has_possession(E, Theme1, Theme2)*. În acest exemplu, *Theme1* este „*nine police officers*” iar *Theme2* este „*a role in the death of 19 people*”; așadar, situația descrisă în text („*their role*”, care poate fi extins la „*the policeman's role*” prin rezoluția anaforei) susține ipoteza și duce la ENTAILMENT, conform pasului **3.a.** din algoritm.

Din cauza modului în care sunt definite clasele Levin, este posibil ca membrii aceleiași clase să nu fie sinonime, ci, mai mult, să fie antonime, așa cum se poate vedea în exemplul 3.

Ex. 3: Antonimie în aceeași clasă VN

T: *BEIJING — China has rejected Coca-Cola Co.'s \$2.3 billion bid to buy a major Chinese juice producer... (China a refuzat oferta de 2.3 miliarde de dolari a Coca-Cola Co. pentru cumpărarea unui mare producător chinez de băuturi răcoritoare)*

H: *Coca-Cola buys Huiyuan Juice Group. (Coca-Cola cumpără Huinan Juice Group)*

Pentru a putea rezolva corect această pereche, verbele *reject* și *buy* din text trebuie relaționate. Verbul *reject* face parte din clasa *approve-77*; cazul întâlnit în **T** este descris sintactic de cadrul *Agent V Proposition* și semantic de *approve(during(E), Agent, Proposition)*. Deoarece verbul în discuție este antonim cu predicatul semantic folosit pentru descrierea sa, descrierea semantică devine *not(approve (during(E), Agent, Proposition))*. Clasa atașată verbului *buy* din **H** este aceeași cu cea din **T**. Pentru rezolvarea corectă a acestui punct, algoritmul trebuie augmentat cu un set de reguli de compunere de verbe, care să conțină reguli de tipul *not(approve(X))→not(X)*. Acest set de reguli poate fi creat relativ ușor, dat numărul de doar 145 de predicate booleene. Compunerea celor două predicate, *reject* și *buy*, duc la rezultatul corect, CONTRADICTION, conform pasului **3.b.** din algoritm.

În cazul în care nu există legătură între verbe pe baza claselor Levin, descrierea semantică a verbelor poate fi folosită pentru a determina o relație semantică, ca în exemplul 4.

Ex. 4: Potrivire pe baza descompunerii semantice

T: *... Fiat, the Italian car company that wants to acquire a stake in Chrysler. ... (Fiat, compania de automobile italiană care vrea să cumpere o parte din Chrysler)*

H: *Fiat wants to gain possession of a stake in Chrysler. (Fiat vrea să câștige posesia unei părți din Chrysler)*

Deși nu există nici un verb în **T** care să fie în aceeași clasă cu *gain (get-13.5.1)*, descrierea semantică a acestei clase, *has_possession(start(E), ?Source, Theme) transfer(during(E), Theme) has_possession(end(E), Agent, Theme) cause(Agent, E)*, este identică cu descrierea semantică pentru clasa *obtain-13.5.2-1*, care conține verbul *acquire* din **T**. Această potrivire, împreună cu potrivirea structurilor de argumente, duc la ENTAILMENT, conform pasului **3.a.** din algoritm.

VerbNet permite stabilirea de relații sintactice între verbe pe baza descrierii semantice a fiecărui cadru, după cum se observă în exemplul 5.

Ex. 5: Stabilirea de relații de consecință

T: *... Indira Gandhi, the former prime minister who was assassinated in 1984. ... (Indira Ghandi, fostul prim ministru care a fost asasinat în 1984.)*

H: *Indira Gandhi died doing aerobatics. (Indira Ghandi a murit făcând acrobații)*

Descrierea semantică pentru clasa *murder-42.1*, din care face parte verbul *assassinate* este *cause(Agent, E) alive(start(E), Patient) not(alive(result(E), Patient))*, iar pentru clasa *disappearance-48.2*, cea a verbului *die*, este *disappear(during(E), Theme)*. Deoarece *disappear* este sinonim cu *die (not alive)*, iar argumentele celor două descrieri semantice se potrivesc (*Theme* și *Patient* sunt similare), se poate deduce că verbul *die* este o consecință a verbului *kill*. Rezultatul este CONTRADICTION, conform pasului **3.b.** din algoritm, deoarece verbul din **H** are un adjunct (*doing aerobatics*) în contradicție cu adjuncții din **T**.

4.2 Perechi rezolvate folosind L și NA

În cazul problemei de rezolvare a inferențelor textuale cu trei soluții posibile (ENTAILMENT, CONTRADICTION, UNKNOWN), una dintre principalele dificultăți este aceea de a determina cazurile UNKNOWN. Conform intuiției pe baza căreia am făcut această analiză, cazurile UNKNOWN apar, de cele mai multe ori, din două motive: una sau mai multe din entitățile din **H** nu apare în **T** (**L**) sau toate entitățile din **H** apar în **T** dar nu sunt conectate de un același predicat (**NA**). Un caz de entități din **H** care nu sunt conectate în **T** este dat în exemplul 6.

Ex. 6: Entități din **H** care nu sunt conectate în **T**

T: *Currently, there is no specific treatment available against dengue fever, which is the most widespread tropical disease after malaria. ... the mosquitoes that transmit dengue ...* (În acest moment nu există un tratament pentru febra dengue, care este cea mai răspândită boală tropicală după malarie...țânțarii care transmit dengue)

H: *Malaria is the most wide-spread disease transmitted by mosquitoes. (Malaria este cea mai răspândită boală transmisă de țânțarii)*

În cadrul textului, cuvântul *malaria* apare doar o singură dată, și nu este conectat cu noțiunea de *mosquitoes*; conform pasului 2.c. din algoritm, acest caz este UNKNOWN.

În urma analizei, toate cele 22 de cazuri de **NA** și cele 43 de cazuri de **L** sunt UNKNOWN, și constituie 65 din totalul de 69 de astfel de cazuri. Celelalte 4 cazuri se datorează faptului că verbele din **H** nu se regăsesc în **T**, după cum se observă în ex. 7.

Ex. 7: Determinarea de UNKNOWN folosind nepotrivirea verbelor

T: *...The iron had to be heated to a precise cherry red color and beaten by the right combination of hammer blows. Mediocre work could hide problems.* (Fierul trebuia încălzit până ajungea la culoarea corectă și trebuia bătut cu o serie de lovituri cu ciocanul speciale. Munca mediocră putea ascunde probleme.)

H: *Titanic sank in 1912. (Titanicul s- scufundat în 1912.)*

Clasa Levin pentru verbul *sink* este *other_cos-45.4*, și are descrierea semantică *state(result(E), Endstate, Patient)*. Singurul verb din **T** similar semantic este *heat*, care face parte din aceeași clasă ca și verbul *sink*, dar nu are ca argument *Titanic* sau ca adjunct *1912*, fapt ce duce la soluția UNKNOWN, conform pasului 3.c. din algoritm.

4.3 Perechi rezolvate folosind AS

După cum am menționat mai sus, **AS** se referă la rezolvarea inferențelor textuale prin metode diferite de compararea de verbe (cel mai frecvent prin compararea de proprietăți ale entităților din **T** și **H**). Dificultatea acestei abordări provine clasificarea atributelor; soluția în acest caz este utilizarea unei ontologii (testele preliminare au demonstrat că atât WordNet (Fellbaum, 1998) cât și SUMO (Niles, Pease, 2001) sunt utile în acest sens). Un exemplu de utilizare a cunoașterii ontologice pentru a putea rezolva inferențe textuale este dat în exemplul 8.

Ex. 8: Rezolvarea inferențelor textuale folosind cunoaștere ontologică

T: *... Mr. Barnes was an offensive lineman in the old American Football League, ...* (Domnul Barnes a fost atacant în vechea liga de fotbal american)

H: Ernie Barnes was an athlete. (Ernie Barnes a fost atlet)

În acest caz, SUMO conține informația că „lineman” este un tip de atlet, fapt ce duce la soluția ENTAILMENT pentru această pereche; de asemenea, verbele copulative trebuie văzute ca o modalitate de asignare de proprietăți.

5. Concluzii

În această lucrare am descris o nouă metodă pentru rezolvarea inferențelor textuale folosind semantică predicțională și structuri de trăsături. Pentru a dovedi fezabilitatea acestei metode am realizat un studiu detaliat a 200 de perechi de inferențe din setul de test din RTE-5. Pentru determinarea semanticii predicționale am examinat utilitatea VerbNet, cu ajutorul căruia pot fi rezolvate direct 38%. De asemenea, am dat și o primă variantă a unui algoritm de rezolvare a inferențelor textuale folosind VerbNet. Un alt avantaj al VN este acela că este aliniat cu o serie de resurse ca FrameNet și WordNet, ceea ce înseamnă că transpunerea sa în limba română poate fi făcută relativ ușor (principalele dificultăți sunt legate de alinierea cadrelor sintactice din română cu cele din engleză). Dacă o variantă de VerbNet pentru limba română, algoritmul propus în această lucrare ar putea rezolva direct perechi de inferențe în limba română.

Întrucât utilitatea VerbNet, în conjuncție cu un context de aplicare (un set de reguli), pentru rezolvarea de inferențe a fost dovedită, iar utilitatea ontologiilor pentru extragerea de seturi de atribute a fost arătată, scopul în acest moment este implementarea acestor idei într-un sistem funcțional. De asemenea, urmează să dezvoltăm o metodă pentru transpunerea VerbNet-ului în limba română pe baza alinierii cu WordNet.

Referințe bibliografice

- Baker, C. F. Fillmore, C. J., Lowe, J. B. (1998). The Berkeley FrameNet project. In *Proceedings of the 17th International Conference on Computational Linguistics (COLING/ACL-98)*, pages 86–90, Montreal.
- Burchardt, A., Frank, A. (2006). Approximating Textual Entailment with LFG and FrameNet Frames. *Proceedings of the 2nd PASCAL Workshop on Recognising Textual Entailment*, Venice, Italy.
- Dagan, I., Glickman, O. (2004). Probabilistic textual entailment: Generic applied modeling of language variability. In *Learning Methods for Text Understanding and Mining*, Grenoble, France.
- Dagan, I. Magnini, B. Glickman, O. (2005). The PASCAL Recognising Textual Entailment Challenge. In *Proceedings of the First PASCAL Challenge Workshop for Recognising Textual Entailment*, pages 1–9, 11–13 April, 2005, Southampton, U.K.
- Fellbaum, C. editor. (1998). WordNet: An Electronic Lexical Database. *Language, Speech and Communications*. MIT Press, Cambridge, Massachusetts.
- Hickl, A., Bensley, J. (2007). A Discourse Commitment-Based Framework for Recognising Textual Entailment. In *Proceedings of the ACL-PASCAL Workshop*

- on Textual Entailment and Paraphrasing*. Pages 185-190. 28-29 June, Prague, Czech Republic.
- Hickl, A. (2008). Using Discourse Commitments to Recognize Textual Entailment. *Proceedings of the 22nd Conference on Computational Linguistics (COLING 2008)*. Manchester, United Kingdom.
- Kipper, K., Dang H. T., Palmer, M. (2000). Class-based construction of a verb lexicon. In *AAAI/IAAI*, pages 691–696.
- Kipper-Schuler, K. (2005). VerbNet: A broad-coverage, comprehensive verb lexicon. *Ph.D. thesis, Computer and Information Science Dept., University of Pennsylvania, Philadelphia, PA, June*
- Kipper-Schuler, K., Korhonen, A., Ryant, N., Palmer, M. (2006). Extending VerbNet with Novel Verb Classes. *Fifth International Conference on Language Resources and Evaluation (LREC 2006)*. Genoa, Italy. June, 2006.
- Korhonen, A., Briscoe, T. (2004). Extended Lexical-Semantic Classification of English Verbs . In *Proceedings of the HLT/NAACL Workshop on Computational Lexical Semantics*, Boston, MA
- Levin, B. (1993). English Verb Classes and Alternations: A Preliminary Investigation, *University of Chicago Press, Chicago, IL*.
- Niles, I., Pease, A. (2001). Towards a standard upper ontology. In *Proceedings of the 2nd International Conference on Formal Ontology in Information Systems FOIS-2001*.
- Moens, M., Steedman, M. (1988). Temporal Ontology and Temporal Reference. *Computational Linguistics*, 14:15–38.
- Palmer, M., Gildea, M., Kingsbury, P. (2005). The Proposition Bank: An Annotated Corpus of Semantic Roles, *Computational Linguistics*, 31:1 , pp. 71-105, March, 2005.

SUMMARIES/ABSTRACTS

METHODOLOGY FOR ESTABLISHING AND ANALYZING ANNOTATED SPEECH CORPUSES – THE CASE OF THE CORPUS SROL

HORIA-NICOLAI TEODORESCU

*Institute for Computer Science of Romanian Academy, Iași, România
Technical University of Iași, Faculty of Electronics, Telecommunications and
Information Technology, Iași – România*

hteodor@etti.tuiasi.ro

Abstract

The purpose of the paper is to thoroughly describe a comprehensive methodology for the establishment and analysis of a spoken language corpus, with reference to the Romanian language. While it is of pragmatic and methodological essence, the paper also aims to draw the attention that previously developed corpora that were used in linguistic inference are unsatisfactorily from the current requirements, hence leading to erroneous conclusions regarding some of the characteristics of spoken languages.

NEW TECHNIQUES FOR THE IDENTIFICATION OF VOWEL AREAS UTTERED IN ROMANIAN LANGUAGE

M. ZBANCIOC^{1,2}, H.N. TEODORESCU^{1,2}, M. FERARU¹

¹*Institute for Computer Science of Romanian Academy, Iași, România*

²*Technical University of Iași, Faculty of Electronics, Telecommunications and
Information Technology, Iași – România*

{hteodor, zmarius}@etti.tuiasi.ro

Abstract

A new set of segmentation techniques is presented, techniques used in order to identify vocalic areas, information that will subsequently be used by the extraction instruments for the fundamental frequency F_0 , moreover for determining the value of the formants F_1, \dots, F_4 . The segmentation phase is important because its errors will have a direct effect on the performances of the pitch extractor. The segmentation precision of our instrument is compared to that of the Praat tool, using high precision annotated files. In order to reduce the running time, a series of optimizations were made to the computation of the autocorrelation function, by applying several recurrent algorithms.

SUMMARIES/ABSTRACTS

METHODOLOGICAL ASPECTS OF ORGANIZING DATA AND STATISTICAL ANALYSIS OF EMOTIONAL VOICE

HORIA-NICOLAI TEODORESCU^{1,2}, IOAN PĂVĂLOI¹, MONICA FERARU¹

¹*Institute for Computer Science of Romanian Academy, Iași, România*

²*Technical University of Iași, Faculty of Electronics, Telecommunications and
Information Technology, Iași – România*

{hteodor}@etti.tuiasi.ro

Abstract

We present a methodology and a program for statistical analysis of emotional voices. Moreover, we present some preliminary results regarding the organization of data in an application for voice signal files. The program allows statistical analysis of the vowels and semivowels of formantic characteristics on the subclasses of data files selected according to user characteristics. The application allows a degree of refinement in the voice sounds analysis.

ROMANIAN INTONATIONAL CONTOUR EDITOR BASED ON FUNCTIONAL PROSODIC FORMS

VASILE APOPEI, DOINA JITCĂ, OTILIA PĂDURARU

*Institute of Computer Science,
Romanian Academy – Iasi Branch*

vapopei@ iit.tuiasi.ro, jdoina@iit.tuiasi.ro

Abstract

The paper presents an intonational tree editor that may be used in Romanian intonation analysis and synthesis. The editor helps the user in spanning the input phrase into multi level prosodic units in order to generate a certain intonational contour in phrase utterance synthesis. In essence the tool offers commands for defining the tree structure and the unit-objects from the nodes and the leaves. The most important attribute in unit-object defining is a functional one at communication act level that assigns it a certain F0 contour pattern (an elementary melodic contour). By using this editor, the understanding of the intonation model is improved. The section 2 presents the intonation model from its functional and hierarchical perspective. The editor functionalities are presented in section 3 and they refer to the intonational tree file management, to the intonational tree editing and to their conversion into speech outputs corresponding to synthesized input phrase utterances.

SUMMARIES/ABSTRACTS

CORPUS FOR GNATHOPHONY: PROTOCOL, METHODOLOGY, ANNOTATION

HORIA-NICOLAI TEODORESCU^{1,2}, ALINA UNTU¹

*¹Technical University of Iasi, Faculty of Electronics,
Telecommunications and Information Technology, Iași - Romania*

*²Institute for Computer Science of the Romanian Academy, Iași -
Romania;*

{hteodor, auntu}@etti.tuiasi.ro

Abstract

We present aspects regarding the extension of a gnathophonic sounds micro-corpus, examples of formantic and temporal analysis applied to the fricative consonants in prosthesis and edentulous cases, as well as preliminary results.

ELECTRONIC DEVICE FOR PARALLEL PREPROCESSING OF THE VOCAL SIGNAL SPECTRUM

HULEA MIRCEA¹, UNTU ALINA²

*¹Technical University „Gheorghe Asachi” of Iasi, Faculty of Automatic Control and
Computer Science, Iași – Romania;*

*²Technical University „Gheorghe Asachi”, Faculty of Electronics, Telecommunications
and Information Technology, Iași – România;*

mhulea@tuiasi.ro, auntu@etti.tuiasi.ro

Abstract

This paper presents a hardware device used for stimulation of an analogue neural network able to perform speech recognition. Considering that spiking neurons are coincidence detectors, the recognition of phonemes should be performed by detection of concurrent events generated by audio signal. Considering that the vowels are characterized by a number of frequencies which are summed during their reception, the system splits the vocal spectrum of the audio signal into frequency channels. The main advantage of this prototype is that it is adapted to generate the trains of impulses needed by a network of spiking neurons to perform detection of singular vowels and vowel combinations. To increase the vowel recognition accuracy the central frequencies of the filters are chosen taking into account the vowels formants. Due to the fact that the device could be used only on vowel detection, as a future research goal we want to test these operation principles in discrimination of the frequency channels activated during consonants reception.

SUMMARIES/ABSTRACTS

CONTINUOUS FLUX ROMANIAN LINGUISTIC RESOURCES

DAN CRISTEA

*Institute of Computer Science, Romanian Academy, Iași – România;
Faculty of Computer Science, „Al. I. Cuza” University, Iași - România*

dcristea@info.uaic.ro

Abstract

The paper raises the issue of a legislative initiative for acquiring large scale language resources. It militates for a large awareness campaign that would bring to stage the importance of storing and preservation for research purpose, in electronic form, of all textual documents which go to print, daily, in Romania. The paper refines the technological steps for this achievement, and stresses legal aspects that should be taken into consideration.

ELECTRONIC COMMUNICATION AND OUR ORTOGRAPHICAL PROBLEMS - WITHOUT SOLUTIONS?

LUCIAN CHIȘU

National Museum of Romanian Literature, Bucharest

lucianchisu@gmail.com

Abstract

This text is a continuation of the author's intervention, published under the title Electronic communication and our ortographical problems (AUSH, series Journalism, 2003). At the time of his first intervention, the digital language in Romanian had a rather transitory aspect. The author considers that, for its optimal usage, two kinds of interventions would be necessary. On one hand, the intervention of technical and system engineers towards creating the software to adapt the keyboard to the Romanian language. On the other hand, the firm attitude of the specialists (linguists, philologists, the cultural elite) towards respecting the ortographical and literary language norms in the case of the digital language, which is about to become the most consistent form of communication between institutions and individuals. The conclusion of the new research, including the 2003-2010 period, is the following: although software exists and Romanian language keyboards are commercialised, and although laws have been passed sanctioning the bad usage of the literary language, the deviations, mistakes and hesitations from the early days are still present in a rather worrying percentage. Official people websites are indicated, some belonging to the Ministry of Culture and Ministry of Education, where ortographical reglementations, far from being put to use, are avoided, in some cases with the express indication to elude diacritics.

SUMMARIES/ABSTRACTS

FACTORS THAT INFLUENCE THE OPTIMAL MOMENT FOR MIGRATING TO THE CORRECT DIACRITICALS IN THE ROMANIAN LANGUAGE

BOGDAN STĂNCESCU

S.C. Moongate Video Production srl, București – România

bogdan@moongate.ro

Abstract

The author estimates that during 2010-2015 most Romanian language content available in digital form will migrate from using the legacy characters with diacritics (using cedilla, per ISO-8859-2) to the proper characters (using comma, per Unicode 3.0). The paper synthesizes the relevant factors that influence the ideal moment for migration, from case to case.

AUTOMATIC EXTRACTION OF MULTILINGUAL CORPUS

TIBERIU BOROȘ, DAN TUFIȘ, ALEXANDRU CEAUȘU

Research Institute for Artificial Intelligence, Romanian Academy, Bucharest, Romania

{tibi,tufis,aceausu}@racai.ro

Abstract

The lack of resources is the primary impediment in technological progress in the automatic processing of most natural languages. The success of statistical approaches in automatic translation for wide spread languages has shown the importance of collecting and preprocessing large volumes of linguistic resources. Parallel corpuses and recently, comparable corpuses, represent the primary resource for cross and multilingual applications. This article presents a set of tools that can be used in automatic extraction of parallel and strongly comparable corpuses.

SUMMARIES/ABSTRACTS

LEXICAL SEMANTICS MODELING AND ROBUST PARSING FOR ROMANIAN, FRENCH, AND GERMAN THESAURI

NECULAI CURTEANU¹, MIHAI ALEX MORUZ^{1,2}, DIANA TRANDABĂȚ^{1,2}

¹*Institute of Computer Science, Romanian Academy, Iași – România;*

²*Faculty of Computer Science, „Al. I. Cuza” University, Iași - România*

{curteanu, mmoruz, dtrandabat}@iit.tuiasi.ro

Abstract

This paper presents a lexical-semantics cross-linguistic analysis of the largest thesauri currently existing for Romanian, French, and German, and a new, robust and portable method for their entry parsing, based on the technique of Segmentation-Cohesion-Dependency (SCD) configurations. The general idea behind parsing a thesaurus or dictionary is to transform a raw text dictionary entry into an indexable linguistic database of lexical-semantics (sense) trees. The SCD parsing configurations are applied successively on a thesaurus entry in order to identify its lexicographic segments (the first SCD configuration), to extract the tree of its senses and subsenses (the second one), and to parse its atomic and non-atomic sense definitions (the third one). Using previous results on **DLR** (The Romanian Thesaurus – new format), the present paper adapts and applies the SCD-based technology to other four large and complex thesauri: **DAR** (The Romanian Thesaurus – old format), **TLF** (Le Trésor de la Langue Française), **DWB** (Deutsches Wörterbuch – GRIMM), and **GWB** (Göthe-Wörterbuch). This experiment, which was illustrated on significantly large parsed entries of the mentioned thesauri, proved the main achievements of the parsing technique based on SCD-configurations: efficiency, robustness, portability. These qualities, when compared with the classical methods for dictionary entry parsing, derive mainly from at least two SCD-specific facts: the sense tree extraction is performed on the sense marker sequences exclusively, while the processes of sense tree extraction and (atomic) sense definition parsing are shown to be completely separable.

THE ACHIEVEMENT OF A ROMANIAN TREEBANK

CENEL-AUGUSTO PEREZ

Faculty of Letters, „Al. I. Cuza” University, Iași - România

Faculty of Computer Science, „Al. I. Cuza” University, Iași - România

cperez@info.uaic.ro

Abstract

We inventorize the outcomes of the achievement of the Romanian language syntactic corpus, stressing on the problems met during the acquisition of this corpus (we will also mention some steps in the creation of the corpus) and on the solutions of these problems.

SUMMARIES/ABSTRACTS

MONITORING NEOLOGISMS IN NEWSPAPERS WITHIN THE NEOROM PROJECT

ANA-MARIA BARBU

„Iorgu Iordan – Al. Rosetti” Institute of Linguistics, Bucharest, Romania

anamaria.barbu@g.unibuc.ro

Abstract

This paper describes the contribution of the Romanian research team within the international project NEOROM. The project aims at building a lexical database consisting in neologisms which occur in newspapers written in all Romance languages since 2004. The paper has three parts. In the first one, we present the monitoring system adopted by the Romanian team, in the second one we describe the NEOROM interface which allows to register the neological words in the database, whereas in the third part we give a short analysis of the words in the database from the viewpoint of their formation.

EMOTIONS IN WORDS: DEVELOPING A MULTILINGUAL RESOURCE

VICTORIA BOBICEV, VICTORIA MAXIM, TATIANA PRODAN,
NATALIA BURCIU, VICTORIA ANGHELUŞ

Technical University of Moldova, Chişinău, Republic of Moldova
vika@rol.md, maxivica@yahoo.com, tatiana.ursulenco@gmail.com,
natusicb@yahoo.com, lazu_vic@yahoo.com

Abstract

In this paper we describe the process of Russian and Romanian WordNet-Affect creation. WordNet-Affect is a lexical resource created on the basis of the Princeton WordNet which contains information about the emotions that the words convey. It is organized in six basic emotions: anger, disgust, fear, joy, sadness, surprise.

We translated the WordNet-Affect synsets into Russian and Romanian and created an aligned English – Romanian – Russian lexical resource. The resource is freely available for research purposes.

SUMMARIES/ABSTRACTS

A TRAINABLE QA SYSTEM FOR ROMANIAN

DAN ȘTEFĂNESCU, RADU ION, ALEXANDRU CEAUȘU, DAN TUFIȘ, ELENA
IRIMIA, VERGINICA BARBU-MITITELU

Research Institute for Artificial Intelligence, Romanian Academy, Bucharest, Romania
{danstef, radu, aceausu, tufis, elena, vergi}@racai.ro

Abstract

This paper presents a Question-Answering system developed at the Research Institute for Artificial Intelligence (RACAI) in a context provided by a national project. The system was objectively evaluated by the organizers of the CLEF QA competition, under the ResPubliQA task for Romanian. We describe the combination of various relevant factors of the system used for the identification of the most relevant paragraphs as answers to natural language questions. The system is available online on the web page of RACAI. It is fully trainable and its functionality is independent of the genre characterizing the training data.

A MACHINE TRANSLATION SYSTEM FOR ROMANIAN AND FRENCH

MIRABELA NAVLEA, AMALIA TODIRAȘCU

LiLPa, Université de Strasbourg, Strasbourg – France, 22 rue René Descartes, BP
80010, 67084 Strasbourg cedex France;
mirabela.navlea@yahoo.fr, todiras@unistra.fr

Abstract

We present an ongoing project aiming at the development of a Machine Translation system for Romanian and French. We adopt a factored phrase-based statistical machine translation method and we combine several linguistic categories (word form, POS tag, morpho-syntactic properties). We use existing parallel corpora and we also build our own parallel, sentence and word aligned corpora. The corpora are tagged, lemmatized and annotated at chunk level. The paper focus on the presentation of aligned corpora. We present a detailed linguistic analysis of word alignment errors and we define a set of repairing rules, used to improve the quality of the word alignment.

SUMMARIES/ABSTRACTS

INTEROPERABLE AND MULTILINGUAL WEB SERVICES

RADU ION, ALEXANDRU CEAUȘU, DAN ȘTEFĂNESCU, DAN TUFIȘ

Research Institute for Artificial Intelligence, Romanian Academy, Bucharest, Romania

{radu, aceausu, danstef, tufis}@racai.ro

Abstract

Interoperability problems of the linguistic tools and resources are major concerns of current research in Natural Language Processing. Most of the new technologies based on (web) Service Oriented Architectures are important steps towards ensuring interoperability. When web services are language independent or are easily adaptable to new languages, interoperability and multilinguality criteria are largely satisfied. Orchestration of different web services is usually done by choosing a common format of the I/O data. In this respect, several platforms for integrating resources and tools developed in various programming languages and running under various operating systems already appeared. In this paper we present some of the linguistic web services of the Research Institute for Artificial Intelligence of the Romanian Academy (ICIA) adapted to the web services platform WebLicht. This platform is an environment where web services can interact because of the fact that their I/O parameters are converted and required to be in TCF (Text Corpus Format).

HYPONIMY PATTERNS FOR ROMANIAN

VERGINICA BARBU MITITELU

Research Institute for Artificial Intelligence, Romanian Academy, Bucharest, Romania

vergi@racai.ro

Abstract

Hyponymy is a lexical-syntactic relation that was not rigorously studied in Romanian linguistics. Yet, it offers engineers a very efficient way of organizing the lexical material useful in the numerous tasks that they develop and that imply natural language processing. We present here two ways of identifying hyponymy patterns in Romanian, the results obtained, and their evaluation. We foresee uses of these patterns in the end of the article.

SUMMARIES/ABSTRACTS

GENERATIVE MECHANISMS OF DERIVATIONAL MORPHOLOGY

MIRCEA PETIC

*Institute of Mathematics and Computer Science, Academy of Sciences of Moldova,
Chisinau, Republic of Moldova*

mirsha@math.md

Abstract

The article studies the problems of a derivative generator development. The starting point is a lexicon that provides storage of derivatives in a lexicon, which contains not only the graphical representation of derivatives but also their constituent morphemes. This allows studying and formulating rules that would generate derivatives taking into account some restrictions.

CREATING A SEMANTIC ROLE PARSER FOR ROMANIAN

DIANA TRANDABĂȚ, DAN CRISTEA

*Institute of Computer Science, Romanian Academy, Iași, Romania
Faculty of Computer Science, „Al. I. Cuza” University Iași, Romania*

{dtrandabat, dcristea}@info.uaic.ro

Abstract

Semantic parsing, by identifying and classifying semantic entities in context, as well as the relations between them, has a great potential for applications such as text summarization, question-answering or machine translation. Thus, by developing a system that automatically annotates semantic for the Romanian language, this paper represents an important intermediate step towards automatic natural language understanding. For the creation of the semantic role parser for Romanian, it was first necessary to develop a training corpus, annotated with semantic roles. This annotated resource of semantic roles for the Romanian language was produced using as starting point a resource developed for English (FrameNet), and applying an automatic transfer method. Subsequently, using a platform for developing supervised semantic role labeling systems (PASRL), several learning algorithms were trained on the developed corpus and the best obtained model was saved. We discuss the results of applying this technique for Romanian, with the conviction that other languages could also benefit from using the same approach.

SUMMARIES/ABSTRACTS

INTERNET SURVEILLANCE: CONSUMERS OPINION ABOUT CERTAIN PRODUCTS OR EVENTS

ADRIAN IFTENE, ALINA-ELENA MIHĂILĂ, GEORGE-ALEXANDRU VLAD,
GETA STANCU

Faculty of Computer Science, „Al. I. Cuza” University of Iasi, Romania

{adiftene, elena.mihaila, george.vlad, geta.stancu}@info.uaic.ro

Abstract

This paper approaches and offers a solution to a modern day concern: relevant information retrieval and surveillance through the internet. The increasing volume of information that exists on web sites, forums or personal web pages, makes the process of searching for information that is relevant for us to become very complex and time consuming. In addition to their success, social networks like Twitter, MySpace, Facebook, Flickr have more and more users with common interests, and their gathered information has an increasing volume. In such social networks starting from a topic, users can express freely their opinions, add links or relevant photos. The system that we built uses information collected from the Internet and offers to users an easier way to find out positive and negative opinions about a topic. Information is searched on web pages (we prefer blogs, forums or users comments) and on Twitter. Identifying and classifying opinions is done by identifying some emotional triggers and by calculating some valences related to the context in which they appear.

USING VERBS FOR DETERMINING TEXTUAL ENTAILMENT

MIHAI ALEX MORUZ^{1,2}

¹*„Al.I.Cuza” University, Faculty of Computer Science, Iasi – Romania*

²*Romanian Academy, Iasi Branch, Institute for Computer Science – Romania*

mmoruz@info.uaic.ro

Abstract

This paper describes a method of using verbs and their argument structure for solving Textual Entailment for English and Romanian. The analysis carried out over a set of entailment pairs from the RTE-5 test set proves that predicational semantics, as defined by Charles Fillmore and Beth Levin, is useful for solving textual entailment, as 38% of the pairs we analyzed are directly solved by using this type of semantics and 11% of the pairs are solved by using verbal analysis.

INDEX DE AUTORI

Apopei Vasile	45
Angheluș Victoria	141
Bobicev Victoria	141
Barbu Ana-Maria	131
Barbu-Mititelu Verginica	153, 185
Boroș Tiberiu	103
Burciu Natalia	141
Ceașu Alexandru	103, 153, 175
Chișu Lucian	81
Cristea Dan	73, 203
Curteanu Neculai	113
Feraru Monica	23, 35
Hulea Mircea	61
Iftene Adrian	213
Ion Radu	153, 175
Irimia Elena	153
Jitcă Doina	45
Maxim Victoria	141
Mihăilă Alina-Elena	213
Moruz Mihai Alex	113, 223
Navlea Mirabela	165
Păduraru Otilia	45
Păvăloi Ioan	35
Perez Cenel-Augusto	123
Petic Mircea	195
Prodan Tatiana	141
Stancu Geta	213
Stăncescu Bogdan	89
Ștefănescu Dan	153, 175
Teodorescu Horia-Nicolai	13, 23, 35, 51
Todirașcu Amalia	165
Trandabăț Diana	113, 203
Tufiș Dan	103, 153, 175
Untu Alina	51, 61
Vlad George-Alexandru	213
Zbancioc Marius-Dan	23